

A Theory of Local Matching SIFT and Beyond

Hossein Mobahi	Stefano Soatto
CSAIL, MIT	CS Dept., UCLA
Cambridge, MA	Los Angeles, CA
hmobahi@csail.mit.edu	soatto@cs.ucla.edu

Abstract

Why has SIFT been so successful? Why its extension, DSP-SIFT, can further improve SIFT? Is there a theory that can explain both? How can such theory benefit real applications? Can it suggest new algorithms with reduced computational complexity or new descriptors with better accuracy for matching? We construct a general theory of local descriptors for visual matching. Our theory relies on concepts in energy minimization and heat diffusion. We show that SIFT and DSP-SIFT approximate the solution the theory suggests. In particular, DSP-SIFT gives a better approximation to the theoretical solution; justifying why DSP-SIFT outperforms SIFT. Using the developed theory, we derive new descriptors that have fewer parameters and are potentially better in handling affine deformations.

1 Introduction

Questions: Why has SIFT been so successful? Why DSP-SIFT [Dong and Soatto, 2015] can further improve SIFT? Is there a theory that can explain both? How can such theory benefit real applications? Can it suggest new algorithms with reduced computational complexity or new descriptors with better accuracy for matching?

Contributions: We construct a general theory of local descriptors for visual matching. Our theory relies on concepts in *energy minimization* and *heat diffusion*. We show that SIFT and DSP-SIFT approximate the solution the theory suggests. In particular, DSP-SIFT gives a better approximation to the theoretical solution; justifying why DSP-SIFT outperforms SIFT. We derive new algorithms based on this theory. Specifically, we present a computationally efficient approximation to DSP-SIFT algorithm [Dong and Soatto, 2015] by replacing the *sampling* procedure in DSP-SIFT with a closed-form approximation that does not need any sampling. This leads to a significantly faster

algorithm compared to DSP-SIFT. In addition, we derive *new descriptors* directly from this theory. The new descriptors have *fewer parameters* as well as the potential of better handling *affine deformations*, compared to SIFT and DSP-SIFT.

2 Contributions

Throughout this text, isotropic multivariate Gaussian kernel and periodic univariate Gaussian are denoted by k and \tilde{k} ,

$$k_\sigma(\mathbf{x}) \triangleq (2\pi\sigma^2)^{-\dim(\mathbf{x})} e^{-\frac{\|\mathbf{x}\|^2}{2\sigma^2}}, \quad \tilde{k}_\sigma(\phi) \triangleq \sum_{k=-\infty}^{\infty} k_\sigma(\phi + 2\pi k). \quad (1)$$

Consider an image f . Given an origin-centered detected key point \mathbf{x} with assigned scale σ and orientation β . The continuous form of a SIFT descriptor can be expressed as [Dong et al., 2015, Vedaldi and Fulkerson, 2010],

$$h_{SIFT}(\beta, \mathbf{x}) \triangleq \int_{\mathcal{X}} \tilde{k}_{\sigma_r}(\beta - \angle \nabla f(\mathbf{y})) k_{\sigma_d}(\mathbf{y} - \mathbf{x}) \|\nabla f(\mathbf{y})\| d\mathbf{y}, \quad (2)$$

where σ_r resembles the size of each orientation bin, e.g. $\frac{2\pi}{8}$ for 8 bins. σ_d determines the spatial support of the descriptor as a function of σ , e.g. $\sigma_d \triangleq 3\sigma$.

By observing that the above descriptor is *pooling* (weighted averaging) across displacement, [Dong et al., 2015] adds domain size pooling to this construction and suggests,

$$h_{DSP}(\beta, \mathbf{x}) \triangleq \int_{\mathcal{S}} \int_{\mathcal{X}} \tilde{k}_{\sigma_r}(\beta - \angle \nabla f(\mathbf{y})) k_{\sigma_d}(\mathbf{y} - \mathbf{x}) \|\nabla f(\mathbf{y})\| d\mathbf{y} k_{\sigma_s}(\sigma_d - \sigma_{d_0}) d\sigma_d, \quad (3)$$

where $\mathcal{S} \triangleq \mathbb{R}$ and σ_{d_0} is a function of key point's scale σ , e.g. $\sigma_{d_0} \triangleq 3\sigma$.

We develop a theory for descriptor construction by returning to the origin of the problem. Specifically we formulate matching as an energy optimization problem. It is known that the resulted cost function is *nonconvex* for a any realistic matching setup [Dong and Soatto, 2015]. Ideally, one would need to *brute-force search* across all possible transformations to find the right match. This is obviously not practical.

Recently a theory of nonconvex optimization by *heat diffusion* has been proposed [Mobahi and Fisher III, 2015a, Mobahi and Fisher III, 2015b]. The theory offers the best (in a certain sense) tractable solution for nonconvex problems. We show that SIFT and DSP-SIFT *approximate* the energy minimization solution that this theory suggests. By leveraging this connection, we present the following contributions.

The domain-size integration (3) is approximated by *numerical sampling* in [Dong et al., 2015], which is slow. Instead, we propose the following two *closed-form approximations* to this integral.

$$h_{DSP}(\beta, \mathbf{x}) \approx \int_{\mathcal{X}} \tilde{k}_{\sigma_r}(\beta - \angle \nabla f(\mathbf{y})) k_{\sigma_d}(\mathbf{y} - \mathbf{x}) \|\nabla f(\mathbf{y})\| \frac{\sigma_d}{\sqrt{\sigma_d^2 + \|\mathbf{x}\|^2 \sigma_s^2}} e^{\frac{\sigma_s^2 (\|\mathbf{x}\|^2 - \mathbf{x}^T \mathbf{y} - \sigma_d^2)^2}{2\sigma_d^2 (\sigma_d^2 + \|\mathbf{x}\|^2 \sigma_s^2)}}. \quad (4)$$

$$h_{DSP}(\beta, \mathbf{x}) \approx \int_{\mathcal{X}} \tilde{k}_{\sigma_r}(\beta - \angle \nabla f(\mathbf{y})) k_{\sigma_d}(\mathbf{y} - \mathbf{x}) \|\nabla f(\mathbf{y})\| \frac{\sigma_d^2 - \sigma_s \mathbf{x}^T \mathbf{y}}{(\sigma_d^2 + \sigma_s^2 \|\mathbf{x}\|^2)^{\frac{3}{2}}} e^{-\frac{\sigma_d^2 \|\mathbf{x} + \mathbf{y}\|^2 + \sigma_s^2 (\mathbf{x}^T \mathbf{y}^\perp)^2}{2\sigma_d^2 (\sigma_d^2 + \sigma_s^2 \|\mathbf{x}\|^2)}}. \quad (5)$$

In addition, through this theory, we propose a *new descriptor*. This descriptor is *exact* in terms of what this theory suggests. In addition, this descriptor is derived from an *affine matching* formulation, hence may better tolerate affine transforms than SIFT and DSP-SIFT¹. Interestingly, despite handling a broader transformation space, it has fewer parameters than DSP-SIFT. Finally, it is *analytical* and does not need any sampling.

$$h_{heat}(\beta, \mathbf{x}) \triangleq \int_{\mathcal{X}} \frac{e^{-\frac{(\mathbf{y}^T \tilde{\nabla} f(\mathbf{y}))^2}{2\sigma_d^2}} w\left(-\frac{1}{2t} \tilde{\nabla}^T f(\mathbf{y}) (\sigma_d^{-2} \mathbf{y} \mathbf{x}^T + \sigma_a^{-2} \mathbf{I}) \tilde{\mathbf{v}}(\beta, \mathbf{y})\right)}{\|\nabla f(\mathbf{y})\|^2 t^3} \times k_{\sqrt{\sigma_d^2 + \sigma_a^2 \|\mathbf{x}\|^2}} \left(\frac{(\nabla f(\mathbf{y}))^T (\mathbf{x} - \mathbf{y})^\perp}{\|\nabla f(\mathbf{y})\|} \right) d\mathbf{y}, \quad (6)$$

where $\tilde{\nabla} f(\mathbf{y}) \triangleq \frac{\nabla f(\mathbf{y})}{\|\nabla f(\mathbf{y})\|}$, $\tilde{\mathbf{v}}(\beta, \mathbf{y}) \triangleq \frac{(\cos(\beta), \sin(\beta))}{\|\nabla f(\mathbf{y})\|}$, $t \triangleq \sqrt{\frac{(\mathbf{x}^T \tilde{\mathbf{v}}(\beta, \mathbf{y}))^2}{2\sigma_d^2} + \frac{1}{2\sigma_a^2 \|\nabla f(\mathbf{y})\|^2}}$, $w(x) \triangleq \sqrt{\pi} e^{x^2} (1 + 2x^2) \operatorname{erfc}(x) - 2x$, and $(a, b)^\perp \triangleq (b, -a)$.

Note that compared to DSP-SIFT, this descriptor *reduces number of parameters* from two (σ_r and σ_s) to one (σ_a).

An illustration of how h_{heat} differs from h_{SIFT} is as follows. Consider a pair of images, namely image 1 and image 2, each consisting of two patches returned by some key point detector. The goal is to establish correspondence between patches using the ℓ_2 distance between normalized descriptors,

$$d(h_1, h_2) \triangleq \int_0^{2\pi} \int_{\mathcal{X}} \left(\frac{h_1(\beta, \mathbf{x})}{\left(\int_0^{2\pi} \int_{\mathcal{X}} h_1^2(\beta^\dagger, \mathbf{x}^\dagger) d\mathbf{x}^\dagger d\beta^\dagger \right)^{\frac{1}{2}}} - \frac{h_2(\beta, \mathbf{x})}{\left(\int_0^{2\pi} \int_{\mathcal{X}} h_2^2(\beta^\dagger, \mathbf{x}^\dagger) d\mathbf{x}^\dagger d\beta^\dagger \right)^{\frac{1}{2}}} \right)^2 d\mathbf{x} d\beta, \quad (7)$$

where $\mathcal{X} \triangleq \mathcal{X}_1 \cap \mathcal{X}_2$. There are two possible matches: $P_A^1 \leftrightarrow P_A^2 \wedge P_B^1 \leftrightarrow P_B^2$ or $P_A^1 \leftrightarrow P_B^2 \wedge P_B^1 \leftrightarrow P_A^2$; obviously only the former is correct. Distance of matches using SIFT are listed in table 2. Note that SIFT descriptor attains lower distance for the wrong match and thus fails, while the heat descriptor finds the correct match. A visualization of SIFT descriptor and heat descriptor are presented in Figures 2 and 3 respectively.

¹SIFT descriptor gains robustness against displacement by pooling across it. DSP-SIFT gains further robustness against scaling by scale pooling. However, none is robust to affine transform, which would require pooling across more parameters.

	Correct: $P_A^1 \leftrightarrow P_A^2 \wedge P_B^1 \leftrightarrow P_B^2$	Wrong: $P_A^1 \leftrightarrow P_B^2 \wedge P_B^1 \leftrightarrow P_A^2$
SIFT	0.20	0.11
Heat	1.02	1.19

Table 1: Table shows total distance between wrongly matched patches and correctly matched patches. Correctly matched patches need to attain lower distance. SIFT fails to do that in this example, while the new descriptor succeeds.

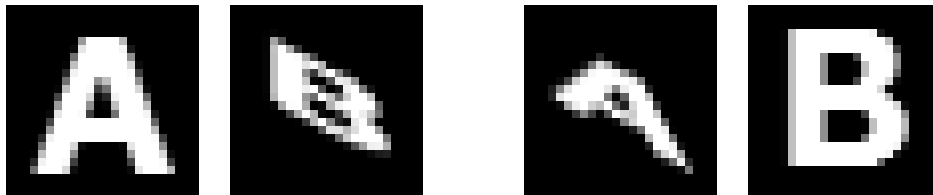


Figure 1: The two patches on the left are considered to belong to image 1, and on the right to image 2.

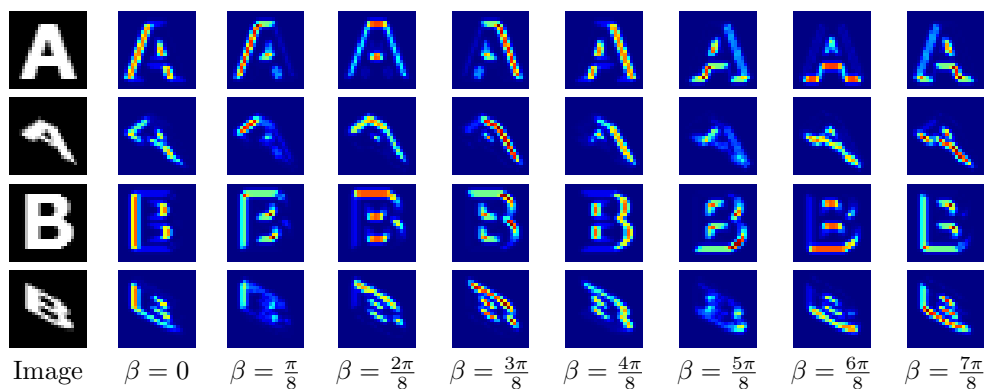


Figure 2: Response map $h_{SIFT}(\beta, \cdot)$ for different choices of β .

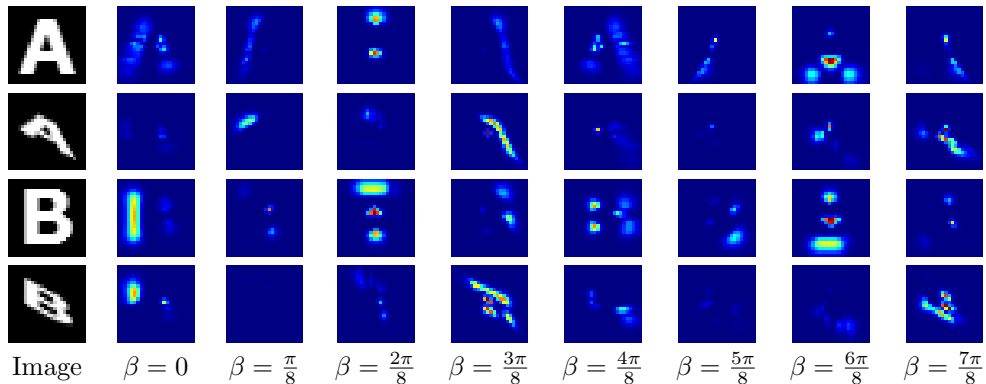


Figure 3: Response map $h_{heat}(\beta, \cdot)$ for different choices of β .

Although this work focuses on SIFT, our diffusion theory can possibly relate and extend other descriptors as well. For example, the recently developed distribution fields [Mears et al., 2013] are similar to (2) and (3), except that instead of histogram of gradient orientation, the histogram of intensity values are used,

$$h_{DF}(l, \mathbf{x}) \triangleq \int_{\mathcal{X}} k_{\sigma_l}(l - f(\mathbf{y})) k_{\sigma_d}(\mathbf{y} - \mathbf{x}) d\mathbf{y}, \quad (8)$$

where σ_l determines the smoothing strength of pixel intensity values. Similar to SIFT arguments, the convolution k_{σ_d} may correspond to diffusion w.r.t. translation, and thus diffusion w.r.t. larger class of transformation, e.g., affine, may lead to geometrically more robust descriptors. Such extensions of distribution fields are not studied in the report, but are subject of future research.

3 Matching as Energy Minimization

For clarity of presentation, we focus on a restricted matching setup with simplifying assumptions. Nevertheless, this setup has enough complexity to make the point on nonconvexity and diffusion.

3.1 Problem Setup

Notation: An image is a map of form $f : \mathcal{X} \rightarrow [0, 1]$, where $\mathcal{X} \subset \mathbb{R}^2$. Similarly, a patch is $p : \mathcal{P} \rightarrow [0, 1]$, where $\mathcal{P} \subseteq \mathcal{X}$, i.e. the map is defined over a subset of the domain \mathcal{X} .

Assumptions: Given a set of patches $p_k : \mathcal{X}_k \rightarrow [0, 1]$ for $k = 1, \dots, n$. We assume that one of these patches, indexed by k^* , appears somewhere in f up to

a geometric transformation $\tau^* : \mathcal{X}_{k^*} \rightarrow \mathcal{X}$ and some reasonable intensity noise²,

$$\exists(k^*, \tau^*) \forall \mathbf{x} \in \mathcal{X}_{k^*} ; f(\tau^*(\mathbf{x})) \approx p_{k^*}(\mathbf{x}) . \quad (9)$$

Objective: The goal is to estimate (k^*, τ^*) . For tractability, the space of τ is *parameterized* by a vector $\boldsymbol{\theta}$. For mathematical convenience, we assume the noise effect is best minimized via ℓ_2 discrepancy,

$$(k^*, \boldsymbol{\theta}^*) \triangleq \underset{(k, \boldsymbol{\theta})}{\operatorname{argmin}} \int_{\mathcal{X}_k} \left(f(\tau(\mathbf{x}; \boldsymbol{\theta})) - p_k(\mathbf{x}) \right)^2 d\mathbf{x} . \quad (10)$$

The tools we later use apply to continuous variables, while (10) involves the integer variable k . However, we can equivalently rewrite the problem in the following continuous form,

$$\begin{aligned} (\mathbf{c}^*, \boldsymbol{\theta}^*) \triangleq \underset{(\mathbf{c}, \boldsymbol{\theta})}{\operatorname{argmin}} \sum_k c_k^2 \int_{\mathcal{X}_k} \left(f(\tau(\mathbf{x}; \boldsymbol{\theta})) - p_k(\mathbf{x}) \right)^2 d\mathbf{x} \\ \text{s.t.} \quad \sum_k c_k = 1 \quad , \quad \forall k ; c_k(1 - c_k) = 0 . \end{aligned} \quad (11)$$

3.2 Intractability

Despite simplicity of the setup, estimation of $(k^*, \boldsymbol{\theta}^*)$ is generally intractable because the optimization problem (11) is *nonconvex*. Hence, local optimization methods may converge to a *local minimum*. In the following, we illustrate this by a toy example. The example involves a univariate signal $f(x)$, a pair of univariate templates $p_1(x)$ and $p_2(x)$ and a translation transform τ so that $f(\tau(x, \theta)) \triangleq f(x - \theta)$. Thus, (11) can be expressed as below, after eliminating c_2 by the equality constraint $c_1 + c_2 = 1$,

$$\begin{aligned} (c_1^*, \theta^*) \triangleq \underset{(c_1, \theta)}{\operatorname{argmin}} \quad c_1^2 \int_{\mathcal{X}_1} (f(x - \theta) - p_1(x))^2 d\mathbf{x} + (1 - c_1)^2 \int_{\mathcal{X}_2} (f(x - \theta) - p_2(x))^2 d\mathbf{x} \\ \text{s.t.} \quad c_1(1 - c_1) = 0 . \end{aligned} \quad (12)$$

The solution c_1^* determines to which template f belongs to; p_1 if $c_1^* = 1$ and p_2 if $c_1^* = 0$.

We proceed by choosing f , p_1 , and p_2 as the blue, green, and red curves in Figure 4-a. Here $\mathcal{X} = [-2, 2]$ and $\mathcal{X}_1 = \mathcal{X}_2 = [-1.2, 1.2]$. The goal is to slide the blue curve to the left or right, such that it coincides with either the green or red curve. Recall from (11) that matching error is examined only over the support of the templates (gray shade). As shown in Figure 4-b, by sliding f to the left by $\theta = 0.25$ units, a perfect match with the green curve is achieved. However, there is no way to attain similar match with the red curve. Thus, by inspection we know that $\mathbf{c}^* = (1, 0)$ and $\theta = 0.25$.

²In this setting each patch p_k may be called a *template*.

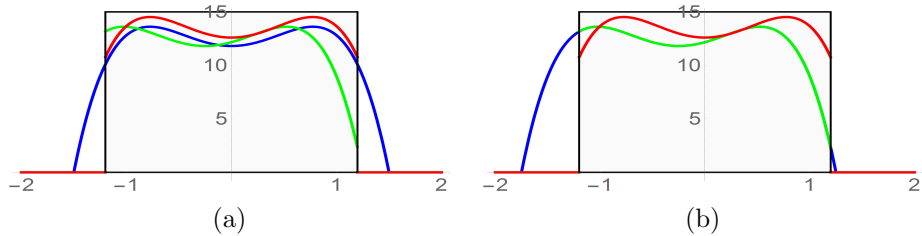


Figure 4: Toy example of signal matching through shift.

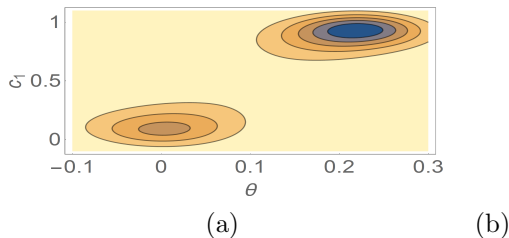


Figure 5: Objective landscape for the toy problem: signal matching through shift.

For visualization purpose, we replace the equality in (12) by a quadratic penalty³. This encompasses both the objective and constraint into a single objective to be visualized. The resulted optimization landscape is shown in Figure 5. A local minimum is apparent around $c_1 = 0, \theta = 0$ while the global minimum is around $c_1 = 1, \theta = 0.25$.

4 Diffusion

One way to approximate the solution of a nonconvex optimization problem is by diffusion and the continuation method. The idea is to *follow* the minimizer of the *diffused* cost function while progressively transforming that function to the original nonconvex cost. It has recently been shown that, this procedure with the choice of the *heat kernel* as the diffusion operator, provides the optimal transformation in a certain sense⁴ [Mobahi and Fisher III, 2015a]. In

³Similar local minima could be obtained for the exact constrained optimization (12) using Lagrange multiplier technique.

⁴It is shown that Gaussian convolution is resulted by the best affine approximation to a nonlinear PDE that generates the convex envelope. Note that computing the convex envelope of a function is generally intractable as well. Thus, it is not surprising that the associated nonlinear PDE lacks a closed form solution. However, by replacing the nonlinear PDE by its best affine approximation, we strike the optimal balance between tractability (closed form solution for the linear PDE) and accuracy of the approximation. The motivation for approximation the convex envelope is that the latter is an optimal object in several senses for the original nonconvex cost function. In particular, global minima of a nonconvex cost are contained in the global minima of its convex envelope.

Algorithm 1 Optimization by Diffusion and Continuation

- 1: Input: $f : \mathcal{X} \rightarrow \mathbb{R}$, Sequence $\sigma_0 > \sigma_1 > \dots > \sigma_n = 0$.
 - 2: $\mathbf{x}_0 =$ global minimizer of $g(\mathbf{x}; \sigma_0)$.
 - 3: **for** $k = 1$ **to** n **do**
 - 4: $\mathbf{x}_k =$ Local minimizer of $g(\mathbf{x}; \sigma_k)$, initialized at \mathbf{x}_{k-1} .
 - 5: **end for**
 - 6: Output: \mathbf{x}_n
-

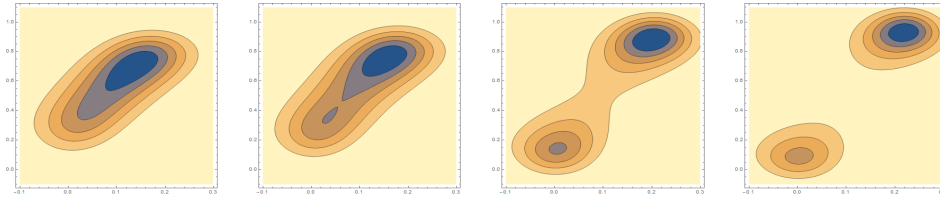


Figure 6: From left to right: diffused cost functions from a large σ toward $\sigma = 0$ for the toy example.

fact, some performance guarantees have been recently developed for this scheme [Mobahi and Fisher III, 2015b]. The procedure is defined more formally below. Given an *unconstrained* and *nonconvex* cost function $h : \mathbb{R}^n \rightarrow \mathbb{R}$ to be minimized. Instead of applying a local optimization algorithm directly to h , we embed h into a family of functions parameterized by σ ,

$$g(\mathbf{x}; \sigma) \triangleq [h \star k_\sigma](\mathbf{x}), \quad (13)$$

where \star is the convolution operator and $k_\sigma(\mathbf{x})$ is the Gaussian function with zero mean and covariance $\sigma^2 \mathbf{I}$. The Gaussian convolution appears here due to the known analytical solution form of the heat diffusion. Observe that $\lim_{\sigma \rightarrow \infty} g(\cdot; \sigma) = h(\mathbf{x})$. Thus by starting from a large σ and shrinking it toward zero, a sequence of cost function converging to h is obtained. The optimization process then follows the path of the minimizer of $g(\cdot; \sigma)$ through this sequence as listed in Algorithm 1.

Now let us revisit the problem (12). Like before, we use a quadratic penalty to obtain an unconstrained approximate to (12). A sequence of diffused landscapes of this problem is shown in Figure 6. Note that the problem becomes convex, with a unique strict minimizer at the large σ . The solution path originated from that point eventually lands at the global minimum in this example.

5 Deriving SIFT via the Diffusion Theory

Instead of pixel intensity as (11) to guide the matching, we switch to orientation of gradient. This change adds limited robustness to illumination changes [Dong and Soatto, 2015]. Nevertheless, the cost function remains nonconvex

and difficult to minimize. Such nonconvex optimization may be treated via diffusion and continuation by the theory of [Mobahi and Fisher III, 2015a]. We show that *SIFT descriptor* emerges as an approximation to this process when τ is a *similarity transformation*, i.e. $\tau(\mathbf{x}; \boldsymbol{\theta}) \triangleq e^s \mathbf{R}_\alpha \mathbf{x} + \mathbf{b}$, where $\boldsymbol{\theta} \triangleq (\alpha, s, \mathbf{b})$.

The approximation comes from two sources. First, the theory of [Mobahi and Fisher III, 2015a] suggests a continuation method by gradually reducing σ while following the path of the minimizer. SIFT provides an approximation to this process by solving the optimization at only one value of σ , i.e. it terminates after the first iteration of the algorithm suggested by [Mobahi and Fisher III, 2015a]. Second, the cost function is diffused only w.r.t. a *subset* of optimization variables (α and \mathbf{b}). This deviates from the theory in [Mobahi and Fisher III, 2015a] that requires diffusion of the cost function in *all variables*, i.e. to use $[cost \star k_\sigma](\mathbf{c}, \boldsymbol{\theta})$.

5.1 Energy Function

Define the density of gradient orientations of image f as below,

$$h(\beta, \mathbf{x}; f) \triangleq \text{III}(\angle \nabla f(\mathbf{x}) - \beta) \|\nabla f(\mathbf{x})\|, \quad (14)$$

where III denotes the *Dirac comb* of period 2π , i.e. $\text{III}(x) \triangleq \sum_{n=-\infty}^{\infty} \delta(x + 2\pi n)$. The Dirac comb accounts for the periodicity of the angle (gradient orientation). Let the dissimilarity between a pair of density functions over a region \mathcal{X} be expressed as the negated dot product,

$$d(f_1, f_2, \mathcal{X}) \triangleq - \int_0^{2\pi} \int_{\mathcal{X}} h(\beta, \mathbf{x}; f_1) h(\beta, \mathbf{x}; f_2) d\beta d\mathbf{x}. \quad (15)$$

In this setting, the problem of template matching is to find $\boldsymbol{\theta}$ such that the gradient orientations of $f(\tau(\mathbf{x}; \boldsymbol{\theta}))$ match that of a template,

$$\begin{aligned} (\mathbf{c}^*, \boldsymbol{\theta}^*) &\triangleq \underset{(\mathbf{c}, \boldsymbol{\theta})}{\operatorname{argmin}} \sum_k c_k d(f \circ \tau_{\boldsymbol{\theta}}, p_k, \mathcal{X}_k) \\ \text{s.t.} \quad &\sum_k c_k = 1 \quad , \quad \forall k; c_k(1 - c_k) = 0. \end{aligned} \quad (16)$$

Replacing the equality constraint by some penalty function q leads to the following *unconstrained* optimization,

$$cost(\mathbf{c}, \boldsymbol{\theta}) \triangleq q(\mathbf{c}) + \sum_k c_k d(f \circ \tau_{\boldsymbol{\theta}}, p_k, \mathcal{X}_k). \quad (17)$$

5.2 Solution

The goal is to tackle the nonconvex problem (17) using the diffusion and continuation theory of [Mobahi and Fisher III, 2015a]. This would give the algorithm listed in Table 5.2-left.

<p>1: Input: $\sigma_0 > \sigma_1 > \dots > \sigma_n = 0$.</p> <p>2: $(\theta_0, \mathbf{c}_0) = \operatorname{argmin}_{(\theta, \mathbf{c})} [\operatorname{cost} \star k_{\sigma_0}](\mathbf{c}, \theta)$.</p> <p>3: for $k = 1$ to n do</p> <p>4: $(\theta_k, \mathbf{c}_k) = \operatorname{Local\ min\ of}$ $[\operatorname{cost} \star k_{\sigma_k}](\mathbf{c}, \theta)$, initialized at $(\theta_{k-1}, \mathbf{c}_{k-1})$.</p> <p>5: end for</p> <p>6: Output: (θ_n, \mathbf{c}_n)</p>	<p>1: Input: $\sigma_0, \Theta \triangleq \cup_j \{(\alpha_j, s_j, \mathbf{b}_j)\}$.</p> <p>2: $(\theta_0, \mathbf{c}_0) = \operatorname{argmin}_{(\mathbf{c}, \alpha, s, \mathbf{b})} [\operatorname{cost}(\mathbf{c}, \cdot, \cdot, \cdot)](\alpha, s, \mathbf{b})$ s.t. $(\alpha, s, \mathbf{b}) \in \Theta$.</p> <p>3: Output: (θ_0, \mathbf{c}_0)</p>	<p>1: Input: $\sigma_0, \Theta \triangleq \cup_j \{(\alpha_j, s_j, \mathbf{b}_j)\}$.</p> <p>2: $(j^*, k^*) = \operatorname{argmax}_{j, k} \int_0^{2\pi} \int_{\mathcal{X}_k^\dagger} h(\beta, \mathbf{x}; p_k) h_{SIFT}(\beta, \mathbf{x}, ; f_j) d\mathbf{x} d\beta$ s.t. $(\alpha, s, \mathbf{b}) \in \Theta$.</p> <p>3: Output: (j^*, k^*)</p>
--	---	---

Table 2: **Left:** Ideal Minimization Strategy based on [Mobahi and Fisher III, 2015a]. **Middle:** Approximation due to a fixed σ and partial diffusion. **Right:** Equivalence with SIFT up to the approximation (20).

SIFT based matching can be derived by *simplifying* this algorithm as described below,

- **Partial Diffusion:** Instead of diffusion w.r.t. all variables $(\alpha, s, \mathbf{b}, \mathbf{c})$, diffuse the energy function (17) partially, i.e. only with respect to (α, \mathbf{b}) .
- **Fixed σ :** Instead of gradual refinement of the energy function by shrinking σ toward zero, stick to a single choice $\sigma = \sigma_0$.
- **Limited Optimization:** Rather than searching the entire parameter space for (α, s, \mathbf{b}) for the optimal solution, restrict to a small candidate set $\Theta \triangleq \cup_{j=1}^J \{(\alpha_j, s_j, \mathbf{b}_j)\}$. This set is generated *outside of the optimization* loop by a keypoint detector⁵. Consequently, Θ does not necessarily contain the optimal parameter as keypoint estimation is done for each image in isolation and thus separately from the full matching problem.

Applying these simplifications yields the algorithm in Table 5.2-middle. The central optimization in this algorithm is the following,

$$\min_{(\alpha, s, \mathbf{b}) \in \Theta} \min_{\mathbf{c}} [[\operatorname{cost}(\mathbf{c}, (\cdot, s, \cdot)) \star k_{\sigma_d}](\mathbf{b}) \star \tilde{k}_{\sigma_r}](\alpha), \quad (18)$$

where we have replaced the joint optimization $\min_{\mathbf{c}, (\alpha, s, \mathbf{b}) \in \Theta}$ by the equivalent nested form $\min_{(\alpha, s, \mathbf{b}) \in \Theta} \min_{\mathbf{c}}$. The *outer* minimization is trivial; it just loops over the candidates and evaluates the resulted cost to pick the best one. Below we only focus on the *inner* optimization.

Assuming the penalty function $q(\mathbf{c})$ accurately enforces the constraint $c_k \in \{0, 1\}$, the inner optimization becomes a *winner take all* problem; the winning patch p_k to match f is the one which minimizes the following cost,

⁵The location and scale of candidate sets are determined by an *interest point detector*, and the orientation angle is set to the *dominant gradient direction*.

$$\begin{aligned}
k^* &= \underset{k}{\operatorname{argmin}} \left[[d(f \circ \tau(\cdot, s, \cdot), p_k, \mathcal{X}_k) \star k_{\sigma_d}](\mathbf{b}) \star \tilde{k}_{\sigma_r} \right](\alpha) \\
&= \underset{k}{\operatorname{argmax}} \int_0^{2\pi} \int_{\mathcal{X}_k} \left([[h(\beta, \mathbf{x}; f \circ \tau(\cdot, s, \cdot)) \star k_{\sigma_d}](\mathbf{b}) \star \tilde{k}_{\sigma_r}](\alpha) \right) h(\beta, \mathbf{x}; p_k) d\beta d\mathbf{x}.
\end{aligned} \tag{19}$$

We doubt that the convolutions in (19) are computationally tractable⁶. Thus, in order to derive a computationally tractable algorithm, we resort to a closed form approximation to the above convolutions. The approximation is stated in the following lemma.

Lemma 1 *The following approximation holds,*

$$\begin{aligned}
& [[h(\beta, \mathbf{x}; f \circ \tau(\cdot, s, \cdot)) \star k_{\sigma_d}](\mathbf{b}) \star \tilde{k}_{\sigma_r}](\alpha) \\
& \approx -e^s \int_{\mathbb{R}^2} \tilde{k}_{\sigma_r}(\angle \nabla f(\mathbf{y}) - \alpha - \beta) \|\nabla f(\mathbf{y})\| k_{\sigma_d}(\mathbf{y} - e^s \mathbf{R}_\alpha \mathbf{x} - \mathbf{b}) d\mathbf{y}.
\end{aligned}$$

Proof See Appendix B for the proof.

Using this lemma, the computationally intractable optimization (18) is replaced by the following tractable approximation,

$$\begin{aligned}
\max_{(\alpha, s, \mathbf{b}) \in \Theta} \max_k e^s \int_0^{2\pi} \int_{\mathcal{X}_k} \int_{\mathbb{R}^2} \tilde{k}_{\sigma_r}(\angle \nabla f(\mathbf{y}) - \alpha - \beta) \|\nabla f(\mathbf{y})\| k_{\sigma_d}(\mathbf{y} - e^s \mathbf{R}_\alpha \mathbf{x} - \mathbf{b}) d\mathbf{y} \\
\times h(\beta, \mathbf{x}; p_k) d\mathbf{x} d\beta,
\end{aligned} \tag{20}$$

In the inner optimization, since (α, s, \mathbf{b}) is fixed (to some $(\alpha_j, s_j, \mathbf{b}_j)$), the image f can be *warped* prior to optimization by $\tau(\cdot; \alpha_j, s_j, \mathbf{b}_j)$. This allows optimization w.r.t. k to be performed for τ being the *identity transform* (because the effect of $(\alpha_j, s_j, \mathbf{b}_j)$ is already taken care of by the warp)⁷, i.e. $(\alpha = 0, s = 0, \mathbf{b} = \mathbf{0})$. Denoting the warped f due to $(\alpha_j, s_j, \mathbf{b}_j)$ by $f_j \triangleq f \circ \tau_{\theta_j}$, the inner optimization simplifies,

$$\max_k \int_0^{2\pi} \int_{\mathcal{X}_k} \underbrace{\int_{\mathbb{R}^2} \tilde{k}_{\sigma_r}(\angle \nabla f_j(\mathbf{y}) - \beta) \|\nabla f_j(\mathbf{y})\| k_{\sigma_d}(\mathbf{y} - \mathbf{x}) d\mathbf{y}}_{h_{SIFT}(\beta, \mathbf{x}; f_j)} h(\beta, \mathbf{x}; p_k) d\mathbf{x} d\beta$$

⁶We will later show in Section 7.2 that by a different *parameterization* of the geometric transform, we can handle a *larger* class, namely the *affine* transform, and yet are able to derive a *closed form* expression for the convolution integrals.

⁷In this section we do not consider the full optimization loop (shrinking σ). However, if we wanted to do so, the idea of 1. Gradual reduction of the blur σ and 2. Warping by the current estimate of the geometric transform in each iteration, would lead to a *Lucas-Kanade* [Lucas and Kanade, 1981] type algorithm. However, the resulted algorithm performs gradient density matching instead of Lucas-Kanade that relies on pixel intensity matching.

Part of the computation involving f_j is *independent of* p_k and can be precomputed. This precomputed result in fact provides a new *representation* for f_j that matches the definition of h_{SIFT} in (2). This gives the algorithm presented in Table 5.2-Right.

6 Deriving DSP-SIFT via the Diffusion Theory

Here we show that the *DSP-SIFT* descriptor also relates to partial diffusion of the cost function. Specifically, this descriptor can be derived by considering the diffusion w.r.t. the transformation parameters (α, s, \mathbf{b}) . Note that this involves more of the optimization variables in diffusion compared to SIFT (which diffuses w.r.t. (α, \mathbf{b})), and thus provides a better approximation to the theory of [Mobahi and Fisher III, 2015a] which suggests the diffusion must be applied to all optimization variables, i.e. to $(\mathbf{c}, \boldsymbol{\theta})$. This improvement in approximation fidelity could be an explanation of why DSP-SIFT works better than SIFT in practice. However, it still misses diffusion of \mathbf{c} .

The derivation is quite similar to that of SIFT in Section 5. The is to minimize the same energy function in (17). However, on top of diffusion w.r.t. variables (α, \mathbf{b}) we add a Gaussian convolution in s . By linearity of convolution, we can just take the diffused energy we obtained in (20) and put it under the Gaussian convolution in s . This leads to the following expression,

$$\max_{(\alpha, s, \mathbf{b}) \in \Theta} \max_k \int_0^{2\pi} \int_{\mathcal{X}_k} \int_{\mathbb{R}^2} \left[(e \cdot \tilde{k}_{\sigma_r}(\angle \nabla f(\mathbf{y}) - \alpha - \beta) \|\nabla f(\mathbf{y})\| k_{\sigma_d}(\mathbf{y} - e \cdot \mathbf{R}_\alpha \mathbf{x} - \mathbf{b}) \star k_{\sigma_s})(s) \right] d\mathbf{y} h(\beta, \mathbf{x}; p_k) d\mathbf{x} d\beta. \quad (22)$$

7 Implication for Future Algorithms

7.1 Closed Form Approximations for Domain Size Pooling

In DSP-SIFT, pooling over the scale is done numerically via sampling [Dong and Soatto, 2015]. Our theory suggests that the scale pooling should also be performed by Gaussian convolution, i.e. (22). Using this form, we present a closed-form approximation, which consequently does *not require any sampling*. Whether or not this approximation provides a satisfactory fidelity must be investigated by experiments.

Recall energy minimization formulation of DSP-SIFT (22) has the following form,

$$\max_{(\alpha, s, \mathbf{b}) \in \Theta} \max_k \int_0^{2\pi} \int_{\mathcal{X}_k} \int_{\mathbb{R}^2} \left([(e \cdot k_{\sigma_d}(\mathbf{y} - e \cdot \mathbf{R}_\alpha \mathbf{x} - \mathbf{b}) \star k_{\sigma_s})(s)] \tilde{k}_{\sigma_r}(\angle \nabla f(\mathbf{y}) - \alpha - \beta) \|\nabla f(\mathbf{y})\| d\mathbf{y} h(\beta, \mathbf{x}; p_k) d\mathbf{x} \right) \quad (23)$$

This convolution does not have a closed form. However, we consider approximating e^s by its linearized form around $s = 0$ (identity scaling transform), i.e. $e^s \approx 1 + s$. Then the convolution will have a closed form. Below we present two approximations based on this idea.

Linearizing only the inner e^s :

Proposition 2

$$\begin{aligned} & [(e \cdot k_{\sigma_d}(\mathbf{y} - (1 + \cdot)\mathbf{R}_\alpha \mathbf{x} - \mathbf{b})) \star k_{\sigma_s}](s) \\ &= k_{\sigma_d}(\mathbf{y} - \mathbf{R}_\alpha \mathbf{x} - \mathbf{b}) \\ & \quad \times k_{\frac{\sigma_d}{\|\mathbf{x}\|}}^{-1} \left(1 + \frac{(\mathbf{R}_\alpha \mathbf{x})^T (\mathbf{b} - \mathbf{y}) - \sigma_d^2}{\|\mathbf{x}\|^2}\right) \times k_{\sqrt{\sigma_s^2 + \frac{\sigma_d^2}{\|\mathbf{x}\|^2}}} \left(s + 1 + \frac{(\mathbf{R}_\alpha \mathbf{x})^T (\mathbf{b} - \mathbf{y}) - \sigma_d^2}{\|\mathbf{x}\|^2}\right). \end{aligned}$$

Proof See Appendix A for the proof.

In particular, when the region is already warped, we can set $(\alpha, s, \mathbf{b}) = (0, 0, \mathbf{0})$. This allows the template matching solution (23) as below,

$$\begin{aligned} & \max_{j,k} \int_0^{2\pi} \int_{\mathcal{X}_k} \int_{\mathbb{R}^2} k_{\sigma_d}(\mathbf{y} - \mathbf{x}) \times k_{\frac{\sigma_d}{\|\mathbf{x}\|}}^{-1} \left(1 - \frac{\mathbf{x}^T \mathbf{y} + \sigma_d^2}{\|\mathbf{x}\|^2}\right) \times k_{\sqrt{\sigma_s^2 + \frac{\sigma_d^2}{\|\mathbf{x}\|^2}}} \left(1 - \frac{\mathbf{x}^T \mathbf{y} + \sigma_d^2}{\|\mathbf{x}\|^2}\right) \\ & \quad \times \tilde{k}_{\sigma_r}(\angle \nabla f_j(\mathbf{y}) - \beta) \|\nabla f_j(\mathbf{y})\| \, d\mathbf{y} h(\beta, \mathbf{x}; p_k) \, d\mathbf{x} \, d\beta. \end{aligned} \quad (24)$$

Linearizing both the inner and outer e^s :

We use the following identity,

$$[(1 + \cdot)k_\sigma(\mathbf{y} + (1 + \cdot)\mathbf{x}) \star k_{scale}](s) \quad (25)$$

$$= \frac{\sigma^2(1+s) - \sigma_{scale} \mathbf{x}^T \mathbf{y}}{2\pi\sigma(\sigma^2 + \sigma_{scale}^2 \|\mathbf{x}\|^2)^{\frac{3}{2}}} e^{-\frac{\sigma^2 \|(1+s)\mathbf{x} + \mathbf{y}\|^2 + \sigma_{scale}^2 (\mathbf{x}^T \mathbf{y})^2}{2\sigma^2(\sigma^2 + \sigma_{scale}^2 \|\mathbf{x}\|^2)}} \quad (26)$$

Thus,

$$[(1 + \cdot)k_{\sigma_d}(\mathbf{y} - \mathbf{b} - (1 + \cdot)\mathbf{R}_\alpha \mathbf{x}) \star k_{\sigma_s}](s) \quad (27)$$

$$= \frac{\sigma_d^2(1+s) + \sigma_s \mathbf{x}^T \mathbf{R}_\alpha^T (\mathbf{y} - \mathbf{b})}{2\pi\sigma_d(\sigma_d^2 + \sigma_s^2 \|\mathbf{x}\|^2)^{\frac{3}{2}}} e^{-\frac{\sigma_d^2 \|(1+s)\mathbf{R}_\alpha \mathbf{x} + \mathbf{y} - \mathbf{b}\|^2 + \sigma_s^2 ((\mathbf{R}_\alpha \mathbf{x})^T (\mathbf{y} - \mathbf{b}))^2}{2\sigma_d^2(\sigma_d^2 + \sigma_s^2 \|\mathbf{x}\|^2)}} \quad (28)$$

In particular, when the region is already warped, we can set $(\alpha, s, \mathbf{b}) = (0, 0, \mathbf{0})$. This allows the template matching solution (23) as below,

$$\begin{aligned} \max_{j,k} \int_0^{2\pi} \int_{\mathcal{X}_k} \int_{\mathbb{R}^2} \frac{\sigma_d^2 + \sigma_s \mathbf{x}^T \mathbf{y}}{(\sigma_d^2 + \sigma_s^2 \|\mathbf{x}\|^2)^{\frac{3}{2}}} e^{-\frac{\sigma_d^2 \|\mathbf{y} - \mathbf{x}\|^2 + \sigma_s^2 (\mathbf{x}^T \mathbf{y}^\perp)^2}{2\sigma_d^2 (\sigma_d^2 + \sigma_s^2 \|\mathbf{x}\|^2)}} & \quad (29) \\ \times \tilde{k}_{\sigma_r}(\angle \nabla f_j(\mathbf{y}) - \beta) \|\nabla f_j(\mathbf{y})\| \, d\mathbf{y} \, h(\beta, \mathbf{x}; p_k) \, d\mathbf{x} \, d\beta & \quad (30) \end{aligned}$$

7.2 Exact Diffusion for Affine Transform

Using the diffusion theory, by using a different *parameterization* for the geometric transformation, we can potentially improved SIFT and DSP-SIFT in two ways.

1. We can extend the descriptor from handling *similarity* transform to *affine* transform.
2. Recall that the computation of the diffusion in SIFT and DSP-SIFT relies on some *approximation*. In addition, regardless of the diffusion theory, DSP-SIFT involves sampling to approximate one of the required *integrals*. The finite sampling process is inaccurate and expensive to compute. Here despite working with a larger transformation space, the new parameterization allows deriving *exact* and *closed form* expression for the diffusion in *all* transformation parameters.

The formulation of the energy function is similar to that of SIFT and DSP-SIFT, except the parameterization. Instead of the similarity transform $\tau(\mathbf{x}; \alpha, s, \mathbf{b}) \triangleq e^s \mathbf{R}_\alpha \mathbf{x} + \mathbf{b}$ we switch to the affine transform $\tau(\mathbf{x}; \mathbf{A}, \mathbf{b}) \triangleq \mathbf{A}\mathbf{x} + \mathbf{b}$, as listed below,

$$\text{cost}(\mathbf{c}, \mathbf{A}, \mathbf{b}) \triangleq q(\mathbf{c}) + \sum_k c_k d(f \circ \tau_{\mathbf{A}, \mathbf{b}, p_k}, \mathcal{X}_k), \quad (31)$$

where the dissimilarity functional d is defined as earlier in (15). Recall that the goal is to tackle the nonconvex problem (31) using the diffusion theory of [Mobahi and Fisher III, 2015a] and similar simplifications as in Section 5, we obtain the following solution for the template matching problem.

$$\begin{aligned} k^* &= \underset{k}{\operatorname{argmin}} \left[[d(f \circ \tau(\dots), p_k, \mathcal{X}_k) \star k_{\sigma_b}](\mathbf{b}) \star \tilde{k}_{\sigma_a}](\mathbf{A}) \right] & (32) \\ &= \underset{k}{\operatorname{argmax}} \int_0^{2\pi} \int_{\mathcal{X}_k} \left([[h(\beta, \mathbf{x}; f \circ \tau(\dots)) \star k_{\sigma_b}](\mathbf{b}) \star \tilde{k}_{\sigma_a}](\mathbf{A}) \right) h(\beta, \mathbf{x}; p_k) \, d\beta \, d\mathbf{x}. \end{aligned}$$

Interesting, we can replace the above convolutions by a *closed form* and *exact* expression. This is stated in the following lemma.

Lemma 3

$$= \frac{[([h(\beta, \mathbf{x}; f \circ \tau(\dots)) \star k_{\sigma_b}](\mathbf{b})) \star k_{\sigma_a}](\mathbf{A})}{e^{-\frac{((\mathbf{b}-\mathbf{y})^T \tilde{\nabla} f(\mathbf{y}))^2}{2\sigma_b^2} - \frac{\|\mathbf{A}^T \tilde{\nabla} f(\mathbf{y})\|^2}{2\sigma_a^2}} w\left(-\frac{\sigma_b^{-2} \tilde{\nabla}^T f(\mathbf{y})(\mathbf{y}-\mathbf{b})\mathbf{x}^T \tilde{\mathbf{v}}(\beta, \mathbf{y}) + \sigma_a^{-2} \tilde{\nabla}^T f(\mathbf{y})\mathbf{A}\tilde{\mathbf{v}}(\beta, \mathbf{y})}{2t}\right)},$$

$$8\sqrt{2}\pi^{\frac{3}{2}}\sigma_b\sigma_a^2\|\nabla f(\mathbf{y})\|^2 t^3,$$

where $\tilde{\nabla} f(\mathbf{y}) \triangleq \frac{\nabla f(\mathbf{y})}{\|\nabla f(\mathbf{y})\|}$, $\tilde{\mathbf{v}}(\beta, \mathbf{y}) \triangleq \frac{\mathbf{v}(\beta)}{\|\nabla f(\mathbf{y})\|}$, and $t \triangleq \sqrt{\frac{(\mathbf{x}^T \tilde{\mathbf{v}}(\beta, \mathbf{y}))^2}{2\sigma_b^2} + \frac{1}{2\sigma_a^2 \|\nabla f(\mathbf{y})\|^2}}$
and $w(x) \triangleq \sqrt{\pi}e^{x^2}(1 + 2x^2) \operatorname{erfc}(x) - 2x$.

Proof See Appendix C for the proof.

Similar to the arguments about SIFT solution in Section 5, the inner optimization in (32) can work with the warped f so that the transformation $\tau(\mathbf{x}; \mathbf{A}, \mathbf{b})$ simplifies to the **identity transform** ($\mathbf{A} = \mathbf{I}, \mathbf{b} = \mathbf{0}$). Letting the warped f be $f_j \triangleq f \circ \tau_{\mathbf{A}_j, \mathbf{b}_j}$, the inner optimization simplifies,

where h_{heat} is defined as the result in lemma 2 (diffused h) with ($\mathbf{A} = \mathbf{I}, \mathbf{b} = \mathbf{0}$), σ_a and σ_b fixed, and all constants dropped,

$$(j^*, k^*) \triangleq \operatorname{argmax}_{j,k} \int_0^{2\pi} \int_{\mathcal{X}_k} h_{heat}(\beta, \mathbf{x}; f_j) h(\beta, \mathbf{x}; p_k) d\mathbf{x} d\beta$$

$$h_{heat}(\beta, \mathbf{x}; f) \triangleq \frac{e^{-\frac{(\mathbf{y}^T \tilde{\nabla} f(\mathbf{y}))^2}{2\sigma_b^2}} w\left(-\frac{1}{2t} \tilde{\nabla}^T f(\mathbf{y}) (\sigma_b^{-2} \mathbf{y}\mathbf{x}^T + \sigma_a^{-2} \mathbf{I}) \tilde{\mathbf{v}}(\beta, \mathbf{y})\right)}{\|\nabla f(\mathbf{y})\|^2 t^3}$$

$$\tilde{\nabla} f(\mathbf{y}) \triangleq \frac{\nabla f(\mathbf{y})}{\|\nabla f(\mathbf{y})\|}$$

$$\tilde{\mathbf{v}}(\beta, \mathbf{y}) \triangleq \frac{\mathbf{v}(\beta)}{\|\nabla f(\mathbf{y})\|}$$

$$t \triangleq \sqrt{\frac{(\mathbf{x}^T \tilde{\mathbf{v}}(\beta, \mathbf{y}))^2}{2\sigma_b^2} + \frac{1}{2\sigma_a^2 \|\nabla f(\mathbf{y})\|^2}}$$

$$w(x) \triangleq \sqrt{\pi}e^{x^2}(1 + 2x^2) \operatorname{erfc}(x) - 2x.$$

References

- [Dong et al., 2015] Dong, J., Karianakis, N., Davis, D., Hernandez, J., Balzer, J., and Soatto, S. (2015). Multi-view feature engineering and learning.
- [Dong and Soatto, 2015] Dong, J. and Soatto, S. (2015). Domain-size pooling in local descriptors: Dsp-sift. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.

- [Lucas and Kanade, 1981] Lucas, B. D. and Kanade, T. (1981). An iterative image registration technique with an application to stereo vision. In *Proceedings of the 7th International Joint Conference on Artificial Intelligence - Volume 2*, IJCAI'81, pages 674–679.
- [Mears et al., 2013] Mears, B., Sevilla-Lara, L., and Learned-Miller, E. G. (2013). Distribution fields with adaptive kernels for large displacement image alignment. In *British Machine Vision Conference, BMVC 2013, Bristol, UK, September 9-13, 2013*.
- [Mobahi and Fisher III, 2015a] Mobahi, H. and Fisher III, J. W. (2015a). On the Link between Gaussian Homotopy Continuation and Convex Envelopes. In Tai, X.-C., Bae, E., Chan, T., and Lysaker, M., editors, *Energy Minimization Methods in Computer Vision and Pattern Recognition*, volume 8932 of *Lecture Notes in Computer Science*, pages 43–56. Springer International Publishing.
- [Mobahi and Fisher III, 2015b] Mobahi, H. and Fisher III, J. W. (2015b). A theoretical analysis of optimization by gaussian continuation. In *Twenty-Ninth AAAI Conference on Artificial Intelligence*.
- [Vedaldi and Fulkerson, 2010] Vedaldi, A. and Fulkerson, B. (2010). Vlfeat: An open and portable library of computer vision algorithms. In *Proceedings of the International Conference on Multimedia*, MM '10, pages 1469–1472. ACM.

Appendix

A Proof of Proposition 1

We proceed with the following identity⁸,

$$\begin{aligned} & e^s \times k_\sigma(\mathbf{y} - (1+s)\mathbf{R}_\alpha \mathbf{x} - \mathbf{b}) \tag{35} \\ = & k_\sigma(\mathbf{y} - \mathbf{R}_\alpha \mathbf{x} - \mathbf{b}) \times k_{\frac{\sigma}{\|\mathbf{R}_\alpha \mathbf{x}\|}}^{-1} \left(1 + \frac{(\mathbf{R}_\alpha \mathbf{x})^T (\mathbf{b} - \mathbf{y}) - \sigma^2}{\|\mathbf{R}_\alpha \mathbf{x}\|^2}\right) \times k_{\frac{\sigma}{\|\mathbf{R}_\alpha \mathbf{x}\|}} \left(s + 1 + \frac{(\mathbf{R}_\alpha \mathbf{x})^T (\mathbf{b} - \mathbf{y}) - \sigma^2}{\|\mathbf{R}_\alpha \mathbf{x}\|^2}\right) \tag{36} \end{aligned}$$

$$= k_\sigma(\mathbf{y} - \mathbf{R}_\alpha \mathbf{x} - \mathbf{b}) \times k_{\frac{\sigma}{\|\mathbf{x}\|}}^{-1} \left(1 + \frac{(\mathbf{R}_\alpha \mathbf{x})^T (\mathbf{b} - \mathbf{y}) - \sigma^2}{\|\mathbf{x}\|^2}\right) \times k_{\frac{\sigma}{\|\mathbf{x}\|}} \left(s + 1 + \frac{(\mathbf{R}_\alpha \mathbf{x})^T (\mathbf{b} - \mathbf{y}) - \sigma^2}{\|\mathbf{x}\|^2}\right) \tag{37}$$

$$\tag{38}$$

In this form, it is now very easy to compute convolution with $k_{\sigma_{scale}}(s)$,

$$\begin{aligned} & [e^\cdot \times k_\sigma(\mathbf{y} - (1 + \cdot)\mathbf{R}_\alpha \mathbf{x} - \mathbf{b}) \star k_{\sigma_{scale}}](s) \tag{39} \\ = & k_\sigma(\mathbf{y} - \mathbf{R}_\alpha \mathbf{x} - \mathbf{b}) \times k_{\frac{\sigma}{\|\mathbf{x}\|}}^{-1} \left(1 + \frac{(\mathbf{R}_\alpha \mathbf{x})^T (\mathbf{b} - \mathbf{y}) - \sigma^2}{\|\mathbf{x}\|^2}\right) \times k_{\sqrt{\sigma_{scale}^2 + \frac{\sigma^2}{\|\mathbf{x}\|^2}}} \left(s + 1 + \frac{(\mathbf{R}_\alpha \mathbf{x})^T (\mathbf{b} - \mathbf{y}) - \sigma^2}{\|\mathbf{x}\|^2}\right) \tag{40} \end{aligned}$$

$$\tag{41}$$

Therefore,

$$\begin{aligned} & \left[\left(e^\cdot \int_{\mathbb{R}^2} \tilde{k}_{\tilde{\sigma}}(\angle \nabla f(\mathbf{y}) - \alpha - \beta) \|\nabla f(\mathbf{y})\| k_\sigma(\mathbf{y} - (1 + \cdot)\mathbf{R}_\alpha \mathbf{x} - \mathbf{b}) d\mathbf{y} \right) \star k_{\sigma_{scale}} \right](s) \tag{42} \\ = & \int_{\mathbb{R}^2} \tilde{k}_{\tilde{\sigma}}(\angle \nabla f(\mathbf{y}) - \alpha - \beta) \|\nabla f(\mathbf{y})\| k_\sigma(\mathbf{y} - \mathbf{R}_\alpha \mathbf{x} - \mathbf{b}) \tag{43} \\ & \times k_{\frac{\sigma}{\|\mathbf{x}\|}}^{-1} \left(1 + \frac{(\mathbf{R}_\alpha \mathbf{x})^T (\mathbf{b} - \mathbf{y}) - \sigma^2}{\|\mathbf{x}\|^2}\right) \times k_{\sqrt{\sigma_{scale}^2 + \frac{\sigma^2}{\|\mathbf{x}\|^2}}} \left(s + 1 + \frac{(\mathbf{R}_\alpha \mathbf{x})^T (\mathbf{b} - \mathbf{y}) - \sigma^2}{\|\mathbf{x}\|^2}\right) d\mathbf{y}. \end{aligned}$$

⁸ We essentially have $e^s \times k_\sigma(\mathbf{y} + (1+s)\mathbf{x})$ which by completing the square of the exponent w.r.t. s can be expressed as below,

$$e^s \times k_\sigma(\mathbf{y} + (1+s)\mathbf{x}) \tag{33}$$

$$= k_\sigma(\mathbf{x} + \mathbf{y}) \times k_{\frac{\sigma}{\|\mathbf{x}\|}}^{-1} \left(1 + \frac{\mathbf{x}^T \mathbf{y} - \sigma^2}{\|\mathbf{x}\|^2}\right) \times k_{\frac{\sigma}{\|\mathbf{x}\|}} \left(s + 1 + \frac{\mathbf{x}^T \mathbf{y} - \sigma^2}{\|\mathbf{x}\|^2}\right), \tag{34}$$

where the first k is 2D, and the next two k 's are 1D.

B Proof of Lemma 1

$$\text{cost}(\mathbf{c}, \boldsymbol{\theta}) \quad (44)$$

$$\triangleq q(\mathbf{c}) + \sum_k c_k d(f \circ \tau_{(\alpha, s, \mathbf{b})}, p_k, \mathcal{X}_k) \quad (45)$$

$$= q(\mathbf{c}) + \sum_k c_k \int_0^{2\pi} \left(h(\beta; f \circ \tau_{(\alpha, s, \mathbf{b})}, \mathcal{X}_k) - h(\beta; p_k, \mathcal{X}_k) \right)^2 d\beta. \quad (46)$$

Note that,

$$h(\beta; f \circ \tau_{(\alpha, s, \mathbf{b})}, \mathcal{X}_k) \quad (47)$$

$$= \int_{\mathcal{X}_k} \mathbb{I}(\angle \nabla(f(e^s \mathbf{R}_\alpha \mathbf{x} + \mathbf{b})) - \beta) \|\nabla f(e^s \mathbf{R}_\alpha \mathbf{x} + \mathbf{b})\| d\mathbf{x} \quad (48)$$

$$= \int_{\mathcal{X}_k} \mathbb{I}(\angle e^s \mathbf{R}_\alpha^T [\nabla f](e^s \mathbf{R}_\alpha \mathbf{x} + \mathbf{b}) - \beta) \|e^s \mathbf{R}_\alpha^T \nabla f(e^s \mathbf{R}_\alpha \mathbf{x} + \mathbf{b})\| d\mathbf{x} \quad (49)$$

$$= \int_{\mathcal{X}_k} \mathbb{I}(\angle \nabla f(e^s \mathbf{R}_\alpha \mathbf{x} + \mathbf{b}) - \alpha - \beta) e^s \|\nabla f(e^s \mathbf{R}_\alpha \mathbf{x} + \mathbf{b})\| d\mathbf{x}, \quad (50)$$

where (49) uses the chain rule of derivative $\nabla(f(\mathbf{A}\mathbf{x} + \mathbf{b})) = \mathbf{A}^T([\nabla f](\mathbf{A}\mathbf{x} + \mathbf{b}))$. Also, for any $a > 0$, (50) uses the identities $\|a\mathbf{R}\mathbf{x}\| = a\|\mathbf{x}\|$ and $\angle a\mathbf{R}\mathbf{x} = \alpha + \angle \mathbf{x}$. Thus, it follows that,

$$\begin{aligned} & \text{cost}(\mathbf{c}, \boldsymbol{\theta}) \quad (51) \\ &= q(\mathbf{c}) + \sum_k c_k \int_0^{2\pi} \left(\int_{\mathcal{X}_k} \mathbb{I}(\angle \nabla f(e^s \mathbf{R}_\alpha \mathbf{x} + \mathbf{b}) - \alpha - \beta) e^s \|\nabla f(e^s \mathbf{R}_\alpha \mathbf{x} + \mathbf{b})\| d\mathbf{x} - h(\beta; p_k, \mathcal{X}_k) \right)^2 d\beta \end{aligned} \quad (52)$$

By the linearity of the convolution operator and the unity of Gaussian's total mass, smoothed cost amounts only to replacing $\left(\int_{\mathcal{X}_k} \mathbb{I}(\angle \nabla f(e^s \mathbf{R}_\alpha \mathbf{x} + \mathbf{b}) - \alpha - \beta) e^s \|\nabla f(e^s \mathbf{R}_\alpha \mathbf{x} + \mathbf{b})\| d\mathbf{x} - h(\beta; p_k, \mathcal{X}_k) \right)^2$ by its smoothed version. Expansion of the quadratic form yields,

$$\begin{aligned} & \left(\int_{\mathcal{X}_k} \mathbb{I}(\angle \nabla f(e^s \mathbf{R}_\alpha \mathbf{x} + \mathbf{b}) - \alpha - \beta) e^s \|\nabla f(e^s \mathbf{R}_\alpha \mathbf{x} + \mathbf{b})\| d\mathbf{x} - h(\beta; p_k, \mathcal{X}_k) \right)^2 \quad (53) \\ &= \left(\int_{\mathcal{X}_k} \mathbb{I}(\angle \nabla f(e^s \mathbf{R}_\alpha \mathbf{x} + \mathbf{b}) - \alpha - \beta) e^s \|\nabla f(e^s \mathbf{R}_\alpha \mathbf{x} + \mathbf{b})\| d\mathbf{x} \right)^2 + h^2(\beta; p_k, \mathcal{X}_k) \quad (54) \\ & \quad - 2 \left(\int_{\mathcal{X}_k} \mathbb{I}(\angle \nabla f(e^s \mathbf{R}_\alpha \mathbf{x} + \mathbf{b}) - \alpha - \beta) e^s \|\nabla f(e^s \mathbf{R}_\alpha \mathbf{x} + \mathbf{b})\| d\mathbf{x} \right) \times h(\beta; p_k, \mathcal{X}_k) \quad (55) \end{aligned}$$

The first term can be rewritten as below,

$$\left(\int_{\mathcal{X}_k} \text{III}(\angle \nabla f(e^s \mathbf{R}_\alpha \mathbf{x} + \mathbf{b}) - \alpha - \beta) e^s \|\nabla f(e^s \mathbf{R}_\alpha \mathbf{x} + \mathbf{b})\| d\mathbf{x} \right)^2 \quad (56)$$

$$= e^{2s} \int_{\mathcal{X}_k} |\{\mathbf{x}_2 \mid \angle \nabla f(e^s \mathbf{R}_\alpha \mathbf{x}_2 + \mathbf{b}) = \angle \nabla f(e^s \mathbf{R}_\alpha \mathbf{x} + \mathbf{b})\}| \quad (57)$$

$$\times \text{III}(\angle \nabla f(e^s \mathbf{R}_\alpha \mathbf{x} + \mathbf{b}) - \alpha - \beta) \|\nabla f(e^s \mathbf{R}_\alpha \mathbf{x} + \mathbf{b})\|^2 d\mathbf{x}. \quad (58)$$

Note that $|\{\mathbf{x}_2 \mid \angle \nabla f(e^s \mathbf{R}_\alpha \mathbf{x}_2 + \mathbf{b}) = \angle \nabla f(e^s \mathbf{R}_\alpha \mathbf{x} + \mathbf{b})\}| \geq 1$ because at least there is one such \mathbf{x}_2 for which $\angle \nabla f(e^s \mathbf{R}_\alpha \mathbf{x}_2 + \mathbf{b}) = \angle \nabla f(e^s \mathbf{R}_\alpha \mathbf{x} + \mathbf{b})$ holds, that is $\mathbf{x}_2 = \mathbf{x}$. However, we assume that the cardinality of the set is exactly one, i.e. besides $\mathbf{x}_2 = \mathbf{x}$, there is no other choice for \mathbf{x}_2 so that the condition $\angle \nabla f(e^s \mathbf{R}_\alpha \mathbf{x}_2 + \mathbf{b}) = \angle \nabla f(e^s \mathbf{R}_\alpha \mathbf{x} + \mathbf{b})$ can hold. The rationale is that the variables are continuous and thus their representation has infinite precision. The odds that the gradient orientation at two different points in the image are *exactly* the same is almost impossible, although they might be very close. With this assumption, the quadratic form simplifies as below,

$$\left(\int_{\mathcal{X}_k} \text{III}(\angle \nabla f(e^s \mathbf{R}_\alpha \mathbf{x} + \mathbf{b}) - \alpha - \beta) e^s \|\nabla f(e^s \mathbf{R}_\alpha \mathbf{x} + \mathbf{b})\| d\mathbf{x} - h(\beta; p_k, \mathcal{X}_k) \right)^2 \quad (59)$$

$$= e^{2s} \int_{\mathcal{X}_k} \text{III}(\angle \nabla f(e^s \mathbf{R}_\alpha \mathbf{x} + \mathbf{b}) - \alpha - \beta) \|\nabla f(e^s \mathbf{R}_\alpha \mathbf{x} + \mathbf{b})\|^2 d\mathbf{x} + h^2(\beta; p_k, \mathcal{X}_k) \quad (60)$$

$$- 2e^s \left(\int_{\mathcal{X}_k} \text{III}(\angle \nabla f(e^s \mathbf{R}_\alpha \mathbf{x} + \mathbf{b}) - \alpha - \beta) \|\nabla f(e^s \mathbf{R}_\alpha \mathbf{x} + \mathbf{b})\| d\mathbf{x} \right) \times h(\beta; p_k, \mathcal{X}_k) \quad (61)$$

$$= e^{2s} \int_{\mathcal{X}_k} \int_{\mathbb{R}^2} \text{III}(\angle \nabla f(\mathbf{y}) - \alpha - \beta) \|\nabla f(\mathbf{y})\|^2 \delta(\mathbf{y} - e^s \mathbf{R}_\alpha \mathbf{x} - \mathbf{b}) d\mathbf{y} d\mathbf{x} + h^2(\beta; p_k, \mathcal{X}_k) \quad (62)$$

$$- 2e^s \left(\int_{\mathcal{X}_k} \int_{\mathbb{R}^2} \text{III}(\angle \nabla f(\mathbf{y}) - \alpha - \beta) \|\nabla f(\mathbf{y})\| \delta(\mathbf{y} - e^s \mathbf{R}_\alpha \mathbf{x} - \mathbf{b}) d\mathbf{y} d\mathbf{x} \right) \times h(\beta; p_k, \mathcal{X}_k) \quad (63)$$

where (62) and (63) use the sifting property of the delta function. The goal is to convolve *cost* with a multivariate Gaussian kernel of covariance $\sigma^2 \mathbf{I}$ in variables jointly in (α, \mathbf{b}) . Due to the diagonal form of the covariance, the convolution can be decoupled to that of α and \mathbf{b} .

We first proceed with smoothing w.r.t. \mathbf{b} . By linearity of the convolution operator and that the Gaussian kernel integrates to one, we obtained the following,

$$\begin{aligned}
& [cost(\mathbf{c}, \alpha, s, \cdot) \star k_\sigma](\mathbf{b}) \\
= & q(\mathbf{c}) + \sum_k c_k \int_0^{2\pi} \left[\left(\int_{\mathcal{X}_k} \text{III}(\angle \nabla f](e^s \mathbf{R}_\alpha \mathbf{x} + \cdot) - \alpha - \beta) e^s \|\nabla f](e^s \mathbf{R}_\alpha \mathbf{x} + \cdot)\| d\mathbf{x} - h(\beta; p_k, \mathcal{X}_k) \right)^2 \star k_\sigma \right. \\
= & q(\mathbf{c}) + \sum_k c_k \int_0^{2\pi} \left(\right. \\
& e^{2s} \int_{\mathcal{X}_k} \int_{\mathbb{R}^2} \text{III}(\angle \nabla f(\mathbf{y}) - \alpha - \beta) \|\nabla f(\mathbf{y})\|^2 k_\sigma(\mathbf{y} - e^s \mathbf{R}_\alpha \mathbf{x} - \mathbf{b}) d\mathbf{y} d\mathbf{x} + h^2(\beta; p_k, \mathcal{X}_k) \\
& \left. - 2e^s \left(\int_{\mathcal{X}_k} \int_{\mathbb{R}^2} \text{III}(\angle \nabla f(\mathbf{y}) - \alpha - \beta) \|\nabla f(\mathbf{y})\| k_\sigma(\mathbf{y} - e^s \mathbf{R}_\alpha \mathbf{x} - \mathbf{b}) d\mathbf{y} d\mathbf{x} \right) \times h(\beta; p_k, \mathcal{X}_k) \right) d\beta.
\end{aligned}$$

We now continue by trying to smooth w.r.t. α .

$$\begin{aligned}
& [[cost(\mathbf{c}, \cdot, s, \cdot) \star k_\sigma] (\mathbf{b}) \star k_{\tilde{\sigma}}] (\alpha) \tag{69} \\
= & q(\mathbf{c}) + \sum_k c_k \int_0^{2\pi} \left(\right. \tag{70} \\
& e^{2s} \int_{\mathcal{X}_k} \int_{\mathbb{R}^2} \|\nabla f(\mathbf{y})\|^2 ([\text{III}(\angle \nabla f(\mathbf{y}) - \cdot - \beta) k_\sigma(\mathbf{y} - e^s \mathbf{R}_\alpha \mathbf{x} - \mathbf{b}) \star k_{\tilde{\sigma}}](\alpha)) d\mathbf{y} d\mathbf{x} + h^2(\beta; p_k, \mathcal{X}_k) \tag{71} \\
& \left. - 2e^s \left(\int_{\mathcal{X}_k} \int_{\mathbb{R}^2} \|\nabla f(\mathbf{y})\| ([\text{III}(\angle \nabla f(\mathbf{y}) - \cdot - \beta) k_\sigma(\mathbf{y} - e^s \mathbf{R}_\alpha \mathbf{x} - \mathbf{b}) \star k_{\tilde{\sigma}}](\alpha)) d\mathbf{y} d\mathbf{x} \right) \times h(\beta; p_k, \mathcal{X}_k) \right) \tag{72}
\end{aligned}$$

Computation of the convolution $\text{III}(\angle \nabla f(\mathbf{y}) - \cdot + \beta) k_\sigma(\mathbf{y} - e^s \mathbf{R}_\alpha \mathbf{x} - \mathbf{b}) \star k_{\tilde{\sigma}}$ is *intractable*. However, it can be approximated by applying the convolution only to the III function. The rationale behind this approximation is that Gaussian convolution affects delta function much more than the Gaussian factor⁹.

$$[\text{III}(\angle \nabla f(\mathbf{y}) - \cdot - \beta) k_\sigma(\mathbf{y} - e^s \mathbf{R}_\alpha \mathbf{x} - \mathbf{b}) \star k_{\tilde{\sigma}}](\alpha) \tag{73}$$

$$\approx \left([\text{III}(\angle \nabla f(\mathbf{y}) - \cdot - \beta) \star k_{\tilde{\sigma}}](\alpha) \right) k_\sigma(\mathbf{y} - e^s \mathbf{R}_\alpha \mathbf{x} - \mathbf{b}) \tag{74}$$

$$= \tilde{k}_{\tilde{\sigma}}(\angle \nabla f(\mathbf{y}) - \alpha - \beta) k_\sigma(\mathbf{y} - e^s \mathbf{R}_\alpha \mathbf{x} - \mathbf{b}). \tag{75}$$

Using this approximation, it follows that,

⁹Gaussian smoothing affects high frequency functions more than low frequency ones; essentially it kills high frequency components, while leaving low frequency components intact.

$$[[\text{cost}(\mathbf{c}, \cdot, s, \cdot) \star k_\sigma] (\mathbf{b}) \star k_{\tilde{\sigma}}] (\alpha) \quad (76)$$

$$\approx q(\mathbf{c}) + \sum_k c_k \int_0^{2\pi} \left(\quad (77)$$

$$e^{2s} \int_{\mathcal{X}_k} \int_{\mathbb{R}^2} \tilde{k}_{\tilde{\sigma}}(\angle \nabla f(\mathbf{y}) - \alpha - \beta) \|\nabla f(\mathbf{y})\|^2 k_\sigma(\mathbf{y} - e^s \mathbf{R}_\alpha \mathbf{x} - \mathbf{b}) d\mathbf{y} d\mathbf{x} + h^2(\beta; p_k, \mathcal{X}_k) \quad (78)$$

$$-2e^s \left(\int_{\mathcal{X}_k} \int_{\mathbb{R}^2} \tilde{k}_{\tilde{\sigma}}(\angle \nabla f(\mathbf{y}) - \alpha - \beta) \|\nabla f(\mathbf{y})\| k_\sigma(\mathbf{y} - e^s \mathbf{R}_\alpha \mathbf{x} - \mathbf{b}) d\mathbf{y} d\mathbf{x} \right) \times h(\beta; p_k, \mathcal{X}_k) \quad (79)$$

C Proof of Lemma 2

1. revert to previous proof with Z , just for the propsotion. Do poper variable replacement in the proof.

Use proposition. Justify why Df having no zero component makes sense.

Mention integration w.r.t y is over $X \cap X_k$. We assume that this integral is zero outside of the domain $R^2 - X$.

¹⁰

Proposition 4

$$\begin{aligned} & \delta(r \mathbf{v}(\beta) - \mathbf{A}^T \nabla f(\mathbf{y})) k_\sigma(\mathbf{y} - \mathbf{A}\mathbf{x} - \mathbf{b}) \\ = & k_\sigma\left(\frac{\mathbf{x}^T \mathbf{z} - (\mathbf{y} - \mathbf{b})^T \nabla f(\mathbf{y})}{\|\nabla f(\mathbf{y})\|}\right) k_{\sqrt{\sigma^2 + \sigma^\dagger^2} \|\mathbf{x}\|^2} \left(\frac{(\nabla f(\mathbf{y}))^T (\mathbf{A}\mathbf{x} + \mathbf{b} - \mathbf{y})^\perp}{\|\nabla f(\mathbf{y})\|}\right) k_{\sigma^\dagger \|\nabla f(\mathbf{y})\|} (\mathbf{z} - \mathbf{A}^T \nabla f(\mathbf{y})) \end{aligned} \quad (80)$$

$$\nabla f(\mathbf{y}) \leftrightarrow \mathbf{g} \quad (82)$$

$$-\mathbf{z} \leftrightarrow \mathbf{c} \quad (83)$$

$$\mathbf{b} - \mathbf{y} \leftrightarrow \mathbf{d} \quad (84)$$

$$(85)$$

Proposition 5

$$[(\delta(\cdot^T \mathbf{g} + \mathbf{c}) k_\sigma(\cdot \mathbf{x} + \mathbf{d})) \star k_{\sigma^\dagger}](\mathbf{A}) \quad (86)$$

$$= k_\sigma\left(\frac{\mathbf{g}^T \mathbf{d} - \mathbf{x}^T \mathbf{c}}{\|\mathbf{g}\|}\right) k_{\sqrt{\sigma^2 + \sigma^\dagger^2} \|\mathbf{x}\|^2} \left(\frac{\mathbf{g}^T (\mathbf{A}\mathbf{x} + \mathbf{d})^\perp}{\|\mathbf{g}\|}\right) k_{\sigma^\dagger \|\mathbf{g}\|} (\mathbf{A}^T \mathbf{g} + \mathbf{c}). \quad (87)$$

Proof We first provide an outline of the proof. The δ function can be replaced by the limit of a Gaussian whose variance tends to zero $\lim_{\epsilon \rightarrow 0} k_\epsilon$. Now k_ϵ and k_{σ^\dagger} form the product of Gaussians. The idea is to write this product as a new single Gaussian in \mathbf{A} (because then we know how to convolve two Gaussians). We do this by replacing the pair with a single exponential whose exponent is trivially the sum of the original exponents. Using *completing the square* method for the joint exponent, the center and covariance of the single Gaussian emerges. There is a problem though; the resulted quadratic form will have a *singular* covariance whose inverse does not exist¹¹. We tackle this problem by the *change of coordinate* system.

¹⁰These conditions can be assumed as granted. The gradient is no where perfectly zero in the image. It is perfectly zero outside of the image f , but that can be taken care of by limiting the integration domain of \mathbf{y} from \mathbb{R}^2 to \mathcal{X} . Having $x_1 = 0$ or $x_2 = 0$ has zero measure, and can be removed from the integration w.r.t. \mathbf{x} without affecting the integration result.

¹¹The inverse of the covariance is required as it directly appears in the definition of the Gaussian.

We begin with the coordinate system transform. Since the covariance of the Gaussian kernel is isotropic, the resulted Gaussian is *radially symmetric*, i.e. $k_\sigma(\mathbf{x}) = k_\sigma(\mathbf{R}\mathbf{x})$ for any rotation matrix \mathbf{R} . Consequently, instead of directly smoothing the above expression, we can rotate the coordinate system, smooth in the latter system, and then invert the rotation to obtain the smoothed function in the original coordinate. In particular, we use the following rotation matrix,

$$\mathbf{R} \triangleq \frac{1}{\|\mathbf{x}\| \|\mathbf{g}\|} \begin{bmatrix} g_2 x_2 \operatorname{sign}(g_1 x_1) & -g_2 |x_1| \operatorname{sign}(g_1) & -x_2 |g_1| \operatorname{sign}(x_1) & |g_1 x_1| \\ -g_1 x_2 \operatorname{sign}(g_2 x_1) & g_1 |x_1| \operatorname{sign}(g_2) & -x_2 |g_2| \operatorname{sign}(x_1) & |g_2 x_1| \\ -g_2 x_1 \operatorname{sign}(g_1 x_2) & -g_2 |x_2| \operatorname{sign}(g_1) & x_1 |g_1| \operatorname{sign}(x_2) & |g_1 x_2| \\ g_1 x_1 \operatorname{sign}(g_2 x_2) & g_1 |x_2| \operatorname{sign}(g_2) & x_1 |g_2| \operatorname{sign}(x_2) & |g_2 x_2| \end{bmatrix}. \quad (88)$$

Due to the assumptions $g_1 \neq 0$, $g_2 \neq 0$, $x_1 \neq 0$, and $x_2 \neq 0$, \mathbf{R} is well-defined. Let $\mathbf{a} \triangleq \operatorname{vec}(\mathbf{A})$, i.e. $\mathbf{a} = (a_{11}, a_{12}, a_{21}, a_{22})$, and let $\mathbf{U} \triangleq \mathbf{R}\mathbf{A}$ and $\mathbf{u} \triangleq \operatorname{vec}(\mathbf{U})$. Changing the coordinate system from \mathbf{A} to \mathbf{U} leads to the following identity,

$$k_\epsilon(\mathbf{A}^T \mathbf{g} + \mathbf{c}) k_\sigma(\mathbf{A} \mathbf{x} + \mathbf{d}) \quad (89)$$

$$= k_\epsilon((\mathbf{R}^T \mathbf{U})^T \mathbf{g} + \mathbf{c}) k_\sigma(\mathbf{R}^T \mathbf{U} \mathbf{x} + \mathbf{d}) \quad (90)$$

$$= \frac{1}{\sqrt{2\pi}} e^{-\frac{\epsilon^2 \|\mathbf{x}\|^2 (\mathbf{d}^T \mathbf{g}^\perp)^2 - \sigma^2 \|\mathbf{g}\|^2 (\|\mathbf{g}\|^2 \|\mathbf{d}\|^2 + 2(-\mathbf{d})^T \mathbf{g} + \mathbf{x}^T \mathbf{c})(\mathbf{x}^T \mathbf{c})}{2\sigma^2 \|\mathbf{g}\|^2 (\sigma^2 \|\mathbf{g}\|^2 + \epsilon^2 \|\mathbf{x}\|^2)}} \quad (91)$$

$$\times e^{\frac{(g_1(-d_2) - g_2(-d_1))^2}{2\sigma^2 \|\mathbf{g}\|^2}} \quad (92)$$

$$\times \frac{1}{\|\mathbf{g}\|} k_{\frac{\epsilon}{\|\mathbf{g}\|}} \left(u_2 - \frac{-x_1 c_2 + x_2 c_1}{\|\mathbf{g}\| \|\mathbf{x}\| \operatorname{sign}(g_2 x_1)} \right) \quad (93)$$

$$\times \frac{1}{\|\mathbf{x}\|} k_{\frac{\sigma}{\|\mathbf{x}\|}} \left(u_3 - \frac{g_1(-d_2) - g_2(-d_1)}{\|\mathbf{g}\| \|\mathbf{x}\| \operatorname{sign}(g_1 x_2)} \right) \quad (94)$$

$$\times \frac{1}{\sqrt{\sigma^2 \|\mathbf{g}\|^2 + \epsilon^2 \|\mathbf{x}\|^2}} k_{\frac{\sigma \epsilon}{\sqrt{\sigma^2 \|\mathbf{g}\|^2 + \epsilon^2 \|\mathbf{x}\|^2}}} \left(u_4 - \frac{\epsilon^2 \|\mathbf{x}\|^2 (-\mathbf{d})^T \mathbf{g} - \sigma^2 \|\mathbf{g}\|^2 \mathbf{x}^T \mathbf{c}}{\|\mathbf{g}\| \|\mathbf{x}\| (\sigma^2 \|\mathbf{g}\|^2 + \epsilon^2 \|\mathbf{x}\|^2) \operatorname{sign}(g_2 x_2)} \right) \quad (95)$$

The value of the coordinate transformation is that we can now write this expression as the product of independent Gaussian kernels. Convolution of this expression with the isotropic kernel $k_{\sigma^\dagger}(\mathbf{u})$ is straightforward,

$$\begin{aligned}
& \delta(\mathbf{z} - \mathbf{A}^T \nabla f(\mathbf{y})) k_\sigma(\mathbf{y} - \mathbf{A}\mathbf{x} - \mathbf{b}) \\
\leftrightarrow & \frac{1}{\sqrt{2\pi}} e^{-\frac{\epsilon^2 \|\mathbf{x}\|^2 ((\mathbf{y}-\mathbf{b})^T \nabla f(\mathbf{y}))^2 - \|\nabla f(\mathbf{y})\|^2 \|\mathbf{y}-\mathbf{b}\|^2 - \sigma^2 \|\nabla f(\mathbf{y})\|^2 (\|\nabla f(\mathbf{y})\|^2 \|\mathbf{y}-\mathbf{b}\|^2 - 2(\mathbf{y}-\mathbf{b})^T \nabla f(\mathbf{y}) - \mathbf{x}^T \mathbf{z})(\mathbf{x}^T \mathbf{z}))}{2\sigma^2 \|\nabla f(\mathbf{y})\|^2 (\sigma^2 \|\nabla f(\mathbf{y})\|^2 + \epsilon^2 \|\mathbf{x}\|^2)}} \\
& \times e^{-\frac{(f_1(\mathbf{y})(y_2-b_2) - f_2(\mathbf{y})(y_1-b_1))^2}{2\sigma^2 \|\nabla f(\mathbf{y})\|^2}} \\
& \times \frac{1}{\|\nabla f(\mathbf{y})\|} k_{\sqrt{\sigma^{\dagger 2} + \frac{\epsilon^2}{\|\nabla f(\mathbf{y})\|^2}}} \left(u_2 - \frac{x_1 z_2 - x_2 z_1}{\|\nabla f(\mathbf{y})\| \|\mathbf{x}\| \text{sign}(f_2(\mathbf{y})x_1)} \right) \\
& \times \frac{1}{\|\mathbf{x}\|} k_{\sqrt{\sigma^{\dagger 2} + \frac{\sigma^2}{\|\mathbf{x}\|^2}}} \left(u_3 - \frac{f_1(\mathbf{y})(y_2-b_2) - f_2(\mathbf{y})(y_1-b_1)}{\|\nabla f(\mathbf{y})\| \|\mathbf{x}\| \text{sign}(f_1(\mathbf{y})x_2)} \right) \\
& \times \frac{1}{\sqrt{\sigma^2 \|\nabla f(\mathbf{y})\|^2 + \epsilon^2 \|\mathbf{x}\|^2}} k_{\sqrt{\sigma^{\dagger 2} + \frac{\sigma^2 \epsilon^2}{\sigma^2 \|\nabla f(\mathbf{y})\|^2 + \epsilon^2 \|\mathbf{x}\|^2}}} \left(u_4 - \frac{\epsilon^2 \|\mathbf{x}\|^2 (\mathbf{y}-\mathbf{b})^T \nabla f(\mathbf{y}) + \sigma^2 \|\nabla f(\mathbf{y})\|^2 \mathbf{x}^T \mathbf{z}}{\|\nabla f(\mathbf{y})\| \|\mathbf{x}\| (\sigma^2 \|\nabla f(\mathbf{y})\|^2 + \epsilon^2 \|\mathbf{x}\|^2)} \text{sign}(f_2(\mathbf{y})x_2)} \right)
\end{aligned}$$

Setting $\epsilon \rightarrow 0$, and given that $\sigma > 0$ and $\|\nabla f(\mathbf{y})\| \neq 0$, it follows that,

$$\delta(\mathbf{z} - \mathbf{A}^T \nabla f(\mathbf{y})) k_\sigma(\mathbf{y} - \mathbf{A}\mathbf{x} - \mathbf{b}) \quad (102)$$

$$\leftrightarrow \frac{1}{\sqrt{2\pi}} e^{-\frac{((\mathbf{b}-\mathbf{y})^T \nabla f(\mathbf{y}) + \mathbf{x}^T \mathbf{z})^2}{2\sigma^2 \|\nabla f(\mathbf{y})\|^2}} \quad (103)$$

$$\times \frac{1}{\|\nabla f(\mathbf{y})\|} k_{\sigma^\dagger} \left(u_2 - \frac{x_1 z_2 - x_2 z_1}{\|\nabla f(\mathbf{y})\| \|\mathbf{x}\| \text{sign}(f_2(\mathbf{y})x_1)} \right) \quad (104)$$

$$\times \frac{1}{\|\mathbf{x}\|} k_{\sqrt{\sigma^{\dagger 2} + \frac{\sigma^2}{\|\mathbf{x}\|^2}}} \left(u_3 - \frac{f_1(\mathbf{y})(y_2-b_2) - f_2(\mathbf{y})(y_1-b_1)}{\|\nabla f(\mathbf{y})\| \|\mathbf{x}\| \text{sign}(f_1(\mathbf{y})x_2)} \right) \quad (105)$$

$$\times \frac{1}{\sigma \|\nabla f(\mathbf{y})\|} k_{\sigma^\dagger} \left(u_4 - \frac{\mathbf{x}^T \mathbf{z}}{\|\nabla f(\mathbf{y})\| \|\mathbf{x}\| \text{sign}(f_2(\mathbf{y})x_2)} \right). \quad (106)$$

By inverting the coordinate system from \mathbf{u} to $(a_{11}, a_{12}, a_{21}, a_{22})$ we obtain,

$$k_\sigma \left(\frac{\mathbf{x}^T \mathbf{z} - (\mathbf{y}-\mathbf{b})^T \nabla f(\mathbf{y})}{\|\nabla f(\mathbf{y})\|} \right) k_{\sqrt{\sigma^2 + \sigma^{\dagger 2} \|\mathbf{x}\|^2}} \left(\frac{(\nabla f(\mathbf{y}))^T (\mathbf{A}\mathbf{x} + \mathbf{b} - \mathbf{y})^\perp}{\|\nabla f(\mathbf{y})\|} \right) k_{\sigma^\dagger \|\nabla f(\mathbf{y})\|} (\mathbf{z} - \mathbf{A}^T \nabla f(\mathbf{y})). \quad (107)$$

As a sanity check, we can see that the above expression becomes the same as the original non-smoothed function when $\sigma^\dagger \rightarrow 0$,

$$\begin{aligned}
& \lim_{\sigma^\dagger \rightarrow 0} k_\sigma \left(\frac{\mathbf{x}^T \mathbf{z} - (\mathbf{y} - \mathbf{b})^T \nabla f(\mathbf{y})}{\|\nabla f(\mathbf{y})\|} \right) k_{\sqrt{\sigma^2 + \sigma^{\dagger 2} \|\mathbf{x}\|^2}} \left(\frac{(\nabla f(\mathbf{y}))^T (\mathbf{A}\mathbf{x} + \mathbf{b} - \mathbf{y})^\perp}{\|\nabla f(\mathbf{y})\|} \right) k_{\sigma^\dagger \|\nabla f(\mathbf{y})\|} (\mathbf{z} - \mathbf{A}^T \nabla f(\mathbf{y})) \\
&= k_\sigma \left(\frac{\mathbf{x}^T \mathbf{z} - (\mathbf{y} - \mathbf{b})^T \nabla f(\mathbf{y})}{\|\nabla f(\mathbf{y})\|} \right) k_\sigma \left(\frac{(\nabla f(\mathbf{y}))^T (\mathbf{A}\mathbf{x} + \mathbf{b} - \mathbf{y})^\perp}{\|\nabla f(\mathbf{y})\|} \right) \delta(\mathbf{z} - \mathbf{A}^T \nabla f(\mathbf{y})) \tag{109}
\end{aligned}$$

$$= k_\sigma \left(\frac{\mathbf{x}^T (\mathbf{A}^T \nabla f(\mathbf{y})) - (\mathbf{y} - \mathbf{b})^T \nabla f(\mathbf{y})}{\|\nabla f(\mathbf{y})\|} \right) k_\sigma \left(\frac{(\nabla f(\mathbf{y}))^T (\mathbf{A}\mathbf{x} + \mathbf{b} - \mathbf{y})^\perp}{\|\nabla f(\mathbf{y})\|} \right) \delta(\mathbf{z} - \mathbf{A}^T \nabla f(\mathbf{y})) \tag{110}$$

$$= k_{\sigma,1} \left(\frac{(\nabla f(\mathbf{y}))^T (\mathbf{A}\mathbf{x} + \mathbf{b} - \mathbf{y})}{\|\nabla f(\mathbf{y})\|} \right) k_{\sigma,1} \left(\frac{(\nabla f(\mathbf{y}))^T (\mathbf{A}\mathbf{x} + \mathbf{b} - \mathbf{y})^\perp}{\|\nabla f(\mathbf{y})\|} \right) \delta(\mathbf{z} - \mathbf{A}^T \nabla f(\mathbf{y})) \tag{111}$$

$$= k_{\sigma,2} \left(\frac{1}{\|\nabla f(\mathbf{y})\|} \left((\nabla f(\mathbf{y}))^T (\mathbf{A}\mathbf{x} + \mathbf{b} - \mathbf{y}), (\nabla f(\mathbf{y}))^T (\mathbf{A}\mathbf{x} + \mathbf{b} - \mathbf{y})^\perp \right) \right) \delta(\mathbf{z} - \mathbf{A}^T \nabla f(\mathbf{y})) \tag{112}$$

$$= k_{\sigma,2} \left(\frac{1}{\|\nabla f(\mathbf{y})\|} \|\nabla f(\mathbf{y})\| (\mathbf{A}\mathbf{x} + \mathbf{b} - \mathbf{y}) \right) \delta(\mathbf{z} - \mathbf{A}^T \nabla f(\mathbf{y})) \tag{113}$$

$$= k_\sigma (\mathbf{A}\mathbf{x} + \mathbf{b} - \mathbf{y}) \delta(\mathbf{z} - \mathbf{A}^T \nabla f(\mathbf{y})) \tag{114}$$

$$\tag{115}$$

□

The goal is to convolve $d(f \circ \tau_{(\mathbf{A}, \mathbf{b})}, p_k, \mathcal{X}_k)$ with the Gaussian kernel. By linearity of the convolution operator we obtain,

$$\begin{aligned}
& \left[\left([d(f \circ \tau_{(\cdot, \cdot)}, p_k, \mathcal{X}_k) \star k_{\sigma_b}](\mathbf{b}) \right) \star k_{\sigma_a} \right](\mathbf{A}) \tag{116} \\
& \triangleq \left[\left(\left[- \int_{\mathcal{X}_k} \int_0^{2\pi} h(\beta, \mathbf{x}; f \circ \tau_{(\cdot, \cdot)}) \times h(\beta, \mathbf{x}; p_k) d\beta d\mathbf{x} \star k_{\sigma_b} \right](\mathbf{b}) \right) \star k_{\sigma_a} \right](\mathbf{A} | d\bar{\mathbf{x}}) \\
& = - \int_{\mathcal{X}_k} \int_0^{2\pi} \left(\left[[h(\beta, \mathbf{x}; f \circ \tau_{(\cdot, \cdot)}) \star k_{\sigma_b}](\mathbf{b}) \right] \star k_{\sigma_a} \right)(\mathbf{A}) \times h(\beta, \mathbf{x}; p_k) d\beta d\bar{\mathbf{x}}
\end{aligned}$$

Thus in the following we focus on $h(\beta, \mathbf{x}; f \circ \tau_{(\cdot, \cdot)}) \star k_{\sigma_b} \star k_{\sigma_a}$. We first manipulate $h(\beta, \mathbf{x}; f \circ \tau_{(\mathbf{A}, \mathbf{b})})$ by applying the chain rule of derivate $\nabla(f(\mathbf{A}\mathbf{x} + \mathbf{b})) = \mathbf{A}^T ([\nabla f](\mathbf{A}\mathbf{x} + \mathbf{b}))$ followed by the sifting property of the delta function,

$$h(\beta, \mathbf{x}; f \circ \tau_{(\mathbf{A}, \mathbf{b})}) \tag{119}$$

$$\triangleq \text{III}(\beta - \angle \nabla(f(\mathbf{A}\mathbf{x} + \mathbf{b})) \| \nabla f(\mathbf{A}\mathbf{x} + \mathbf{b}) \|) \tag{120}$$

$$= \text{III}(\beta - \angle \mathbf{A}^T [\nabla f](\mathbf{A}\mathbf{x} + \mathbf{b}) \| \mathbf{A}^T \nabla f(\mathbf{A}\mathbf{x} + \mathbf{b}) \|) \tag{121}$$

$$= \int_{\mathbb{R}^2} \text{III}(\beta - \angle \mathbf{A}^T \nabla f(\mathbf{y}) \| \mathbf{A}^T \nabla f(\mathbf{y}) \|) \delta(\mathbf{y} - \mathbf{A}\mathbf{x} - \mathbf{b}) d\mathbf{y}. \tag{122}$$

Computing the inner convolution, i.e. w.r.t. \mathbf{b} , is straightforward,

$$\begin{aligned}
& [h(\beta, \mathbf{x}; f \circ \tau_{(\mathbf{A}, \cdot)}) \star k_{\sigma_b}](\mathbf{b}) & (123) \\
= & \left[\left(\int_{\mathbb{R}^2} \text{III}(\beta - \angle \mathbf{A}^T \nabla f(\mathbf{y})) \|\mathbf{A}^T \nabla f(\mathbf{y})\| \delta(\mathbf{y} - \mathbf{A}\mathbf{x} - \mathbf{b}) d\mathbf{y} \right) \star k_{\sigma_b} \right](\mathbf{b}) & (124) \\
= & \int_{\mathbb{R}^2} \text{III}(\beta - \angle \mathbf{A}^T \nabla f(\mathbf{y})) \|\mathbf{A}^T \nabla f(\mathbf{y})\| k_{\sigma_b}(\mathbf{y} - \mathbf{A}\mathbf{x} - \mathbf{b}) d\mathbf{y}. & (125)
\end{aligned}$$

The latter can be expressed by the sifting property of the delta function as below,

$$\begin{aligned}
& [h(\beta, \mathbf{x}; f \circ \tau_{(\mathbf{A}, \cdot)}) \star k_{\sigma_b}](\mathbf{b}) & (126) \\
= & \int_{\mathbb{R}^2} \int_{\mathbb{R}^2} \delta(\mathbf{z} - \mathbf{A}^T \nabla f(\mathbf{y})) \text{III}(\beta - \angle \mathbf{z}) \|\mathbf{z}\| k_{\sigma_b}(\mathbf{y} - \mathbf{A}\mathbf{x} - \mathbf{b}) d\mathbf{z} d\mathbf{y} & (127)
\end{aligned}$$

We now apply a change of variable to move from the Cartesian coordinate (z_1, z_2) to the *polar* coordinate (r, ϕ) such that $(z_1, z_2) = (r \cos(\phi), r \sin(\phi))$. This results in replacing $\int_{\mathbb{R}^2} f(z_1, z_2) dz_1 dz_2$ by $\int_0^\infty \int_0^{2\pi} r f(r \mathbf{v}(\phi)) d\phi dr$, where $\mathbf{v}(\phi) \triangleq (\cos(\phi), \sin(\phi))$.

$$\begin{aligned}
& [h(\beta, \mathbf{x}; f \circ \tau_{(\mathbf{A}, \cdot)}) \star k_{\sigma_b}](\mathbf{b}) & (128) \\
= & \int_{\mathbb{R}^2} \int_0^\infty \int_0^{2\pi} r \delta(r \mathbf{v}(\phi) - \mathbf{A}^T \nabla f(\mathbf{y})) \text{III}(\beta - \angle r \mathbf{v}(\phi)) \|r \mathbf{v}(\phi)\| k_{\sigma_b}(\mathbf{y} - \mathbf{A}\mathbf{x} - \mathbf{b}) d\phi dr d\mathbf{y} & (129) \\
= & \int_{\mathbb{R}^2} \int_0^\infty \int_0^{2\pi} r \delta(r \mathbf{v}(\phi) - \mathbf{A}^T \nabla f(\mathbf{y})) \text{III}(\beta - \angle \mathbf{v}(\phi)) r k_{\sigma_b}(\mathbf{y} - \mathbf{A}\mathbf{x} - \mathbf{b}) d\phi dr d\mathbf{y} & (130) \\
= & \int_{\mathbb{R}^2} \int_0^\infty r^2 \delta(r \mathbf{v}(\beta) - \mathbf{A}^T \nabla f(\mathbf{y})) k_{\sigma_b}(\mathbf{y} - \mathbf{A}\mathbf{x} - \mathbf{b}) dr d\mathbf{y}. & (131)
\end{aligned}$$

We are now ready to smooth this form w.r.t. \mathbf{A} . That is, we want to compute convolution of this expression with a multivariate Gaussian in $(a_{11}, a_{12}, a_{21}, a_{22})$ of covariance $\sigma^2 \mathbf{I}$.



Using this result, we can continue as below,

$$\begin{aligned}
& [[\text{cost}(\mathbf{c}, \cdot, \cdot) \star k_\sigma](\mathbf{b}) \star k_{\sigma^\dagger}](\mathbf{A}) \\
= & q(\mathbf{c}) + \sum_k c_k \int_0^{2\pi} \left(\int_{\mathcal{X}_k} \int_{\mathbb{R}^2} \int_{\mathbb{R}^2} \delta(\mathbf{v}^T(\beta)\mathbf{z}) \|\mathbf{z}\|^2 [(\delta(\mathbf{z} - \cdot^T \nabla f(\mathbf{y})) k_\sigma(\mathbf{y} - \cdot \mathbf{x} - \mathbf{b})) \star k_{\sigma^\dagger}](\mathbf{A}) dz d\mathbf{y} d\mathbf{x} + h^2(\beta; p_k, \mathcal{X}_k) \right. \\
& \left. - 2 \left(\int_{\mathcal{X}_k} \int_{\mathbb{R}^2} \int_{\mathbb{R}^2} \delta(\mathbf{v}^T(\beta)\mathbf{z}) \|\mathbf{z}\| [(\delta(\mathbf{z} - \cdot^T \nabla f(\mathbf{y})) k_\sigma(\mathbf{y} - \cdot \mathbf{x} - \mathbf{b})) \star k_{\sigma^\dagger}](\mathbf{A}) dz d\mathbf{y} d\mathbf{x} \right) \times h(\beta; p_k, \mathcal{X}_k) \right) \\
= & q(\mathbf{c}) + \sum_k c_k \int_0^{2\pi} \left(\int_{\mathcal{X}_k} \int_{\mathbb{R}^2} \int_{\mathbb{R}^2} \delta(\mathbf{v}^T(\beta)\mathbf{z}) \|\mathbf{z}\|^2 k_\sigma \left(\frac{\mathbf{x}^T \mathbf{z} - (\mathbf{y} - \mathbf{b})^T \nabla f(\mathbf{y})}{\|\nabla f(\mathbf{y})\|} \right) k_{\sigma^\dagger \|\nabla f(\mathbf{y})\|} (\mathbf{z} - \mathbf{A}^T \nabla f(\mathbf{y})) dz \right. \\
& \times k_{\sqrt{\sigma^2 + \sigma^{\dagger 2} \|\mathbf{x}\|^2}} \left(\frac{(\nabla f(\mathbf{y}))^T (\mathbf{A}\mathbf{x} + \mathbf{b} - \mathbf{y})^\perp}{\|\nabla f(\mathbf{y})\|} \right) d\mathbf{y} d\mathbf{x} + h^2(\beta; p_k, \mathcal{X}_k) \\
& \left. - 2 \left(\int_{\mathcal{X}_k} \int_{\mathbb{R}^2} \int_{\mathbb{R}^2} \delta(\mathbf{v}^T(\beta)\mathbf{z}) \|\mathbf{z}\| k_\sigma \left(\frac{\mathbf{x}^T \mathbf{z} - (\mathbf{y} - \mathbf{b})^T \nabla f(\mathbf{y})}{\|\nabla f(\mathbf{y})\|} \right) k_{\sigma^\dagger \|\nabla f(\mathbf{y})\|} (\mathbf{z} - \mathbf{A}^T \nabla f(\mathbf{y})) dz \right. \right. \\
& \left. \left. \times k_{\sqrt{\sigma^2 + \sigma^{\dagger 2} \|\mathbf{x}\|^2}} \left(\frac{(\nabla f(\mathbf{y}))^T (\mathbf{A}\mathbf{x} + \mathbf{b} - \mathbf{y})^\perp}{\|\nabla f(\mathbf{y})\|} \right) d\mathbf{y} d\mathbf{x} \right) \times h(\beta; p_k, \mathcal{X}_k) \right) d\beta.
\end{aligned}$$

We now apply a change of variable to move from the Cartesian coordinate (z_1, z_2) to the **polar** coordinate (r, ϕ) such that $(z_1, z_2) = (r \cos(\phi), r \sin(\phi))$. This transforms the form $\int_{\mathbb{R}^2} f(z_1, z_2) dz_1 dz_2$ to $\int_0^\infty \int_0^{2\pi} r f(r \cos(\phi), r \sin(\phi)) d\phi dr$.

$$\int_{\mathbb{R}^2} \text{III}(\beta - \angle \mathbf{z}) \|\mathbf{z}\| k_\sigma \left(\frac{\mathbf{x}^T \mathbf{z} - (\mathbf{y} - \mathbf{b})^T \nabla f(\mathbf{y})}{\|\nabla f(\mathbf{y})\|} \right) k_{\sigma^\dagger \|\nabla f(\mathbf{y})\|} (\mathbf{z} - \mathbf{A}^T \nabla f(\mathbf{y})) dz \quad (141)$$

$$= \int_0^\infty \int_0^{2\pi} r \text{III}(\beta - \phi) r k_\sigma \left(\frac{r \mathbf{x}^T \mathbf{v}(\phi) - (\mathbf{y} - \mathbf{b})^T \nabla f(\mathbf{y})}{\|\nabla f(\mathbf{y})\|} \right) k_{\sigma^\dagger \|\nabla f(\mathbf{y})\|} (r \mathbf{v}(\phi) - \mathbf{A}^T \nabla f(\mathbf{y})) dr d\phi \quad (142)$$

$$= \int_0^\infty r^2 k_\sigma \left(r \frac{\mathbf{x}^T \mathbf{v}(\beta)}{\|\nabla f(\mathbf{y})\|} + \frac{(\mathbf{b} - \mathbf{y})^T \nabla f(\mathbf{y})}{\|\nabla f(\mathbf{y})\|} \right) k_{\sigma^\dagger \|\nabla f(\mathbf{y})\|} (r \mathbf{v}(\beta) - \mathbf{A}^T \nabla f(\mathbf{y})) dr \quad (143)$$

$$= \frac{e^{-\frac{((\mathbf{b} - \mathbf{y})^T \nabla f(\mathbf{y}))^2}{2\sigma^2} - \frac{\|\mathbf{A}^T \nabla f(\mathbf{y})\|^2}{2\sigma^{\dagger 2}}} w\left(-\frac{\sigma^{-2} \tilde{\nabla}^T f(\mathbf{y})(\mathbf{y} - \mathbf{b}) \mathbf{x}^T \tilde{\mathbf{v}}(\beta, \mathbf{y}) + \sigma^{\dagger - 2} \tilde{\nabla}^T f(\mathbf{y}) \mathbf{A} \tilde{\mathbf{v}}(\beta, \mathbf{y})}{2t}\right)}{8\sqrt{2}\pi^{\frac{3}{2}} \sigma \sigma^{\dagger 2} \|\nabla f(\mathbf{y})\|^2 t^3}, \quad (144)$$

where $\tilde{\nabla} f(\mathbf{y}) \triangleq \frac{\nabla f(\mathbf{y})}{\|\nabla f(\mathbf{y})\|}$, $\tilde{\mathbf{v}}(\beta, \mathbf{y}) \triangleq \frac{\mathbf{v}(\beta)}{\|\nabla f(\mathbf{y})\|}$, and $t \triangleq \sqrt{\frac{(\mathbf{x}^T \tilde{\mathbf{v}}(\beta, \mathbf{y}))^2}{2\sigma^2} + \frac{1}{2\sigma^{\dagger 2} \|\nabla f(\mathbf{y})\|^2}}$ and $w(x) \triangleq \sqrt{\pi} e^{x^2} (1 + 2x^2) \text{erfc}(x) - 2x$. In (144) we use an elementary identity¹².

¹² We use the identity,

$$\int_0^\infty r^2 k_{\sigma_1,1}(rc_1 + c_2)k_{\sigma_2,2}(rc_3 + c_4) dr = \frac{e^{-\frac{c_2^2}{2\sigma_1^2} - \frac{\|c_4\|^2}{2\sigma_2^2}} (\sqrt{\pi}(1 + 2t_2^2)e^{t_2^2} \operatorname{erfc}(t_2) - 2t_2)}{8\sqrt{2}\pi^{\frac{3}{2}}\sigma_1\sigma_2^2 t_1^3}, \quad (145)$$

for $t_1 \triangleq \sqrt{\frac{c_1^2}{2\sigma_1^2} + \frac{\|c_3\|^2}{2\sigma_2^2}}$ and $t_2 \triangleq \frac{\frac{c_1 c_2}{s_1^2} + \frac{c_3^T c_4}{\sigma_2^2}}{2t_1}$. This identity is derived in two steps:

1. *Completing the square* of the exponent in the integrand.

$$-\frac{(rc_1 + c_2)^2}{2\sigma_1^2} - \frac{\|rc_3 + c_4\|^2}{2\sigma_2^2} = -\frac{1}{2}\left(r + \frac{c_3^T c_4 \sigma_1^2 + c_1 c_2 \sigma_2^2}{\|c_3\|^2 \sigma_1^2 + c_1^2 \sigma_2^2}\right)^2 \left(\frac{c_1^2}{\sigma_1^2} + \frac{\|c_3\|^2}{\sigma_2^2}\right) + \frac{1}{2}\left(\frac{(c_1 c_2 \sigma_2^2 + \sigma_1^2 c_3^T c_4)^2}{\|c_3\|^2 \sigma_1^4 \sigma_2^2 + c_1^2 \sigma_1^2 \sigma_2^4} - \frac{c_2^2}{\sigma_1^2} - \frac{\|c_4\|^2}{\sigma_2^2}\right). \quad (146)$$

2. Using the identity about *Gaussian moments*,

$$\int_0^\infty r^2 e^{-\frac{(r-a_1)^2}{2a_2^2}} dr = a_1 a_2^2 e^{-\frac{a_1^2}{2a_2^2}} + \sqrt{\frac{\pi}{2}} a_2 (a_1^2 + a_2^2) \left(1 + \operatorname{erf}\left(\frac{a_1}{\sqrt{2}a_2}\right)\right). \quad (147)$$