

Learning theory

Lecture 9

David Sontag
New York University

Slides adapted from Carlos Guestrin & Luke Zettlemoyer

Introduction to probability: events

- An **event** is a subset of the outcome space, e.g.

$$E = \{ \text{die with 2, 4, 6} , \text{die with 1, 3, 5} , \text{die with 2, 4, 6} \} \quad \text{Even die tosses}$$

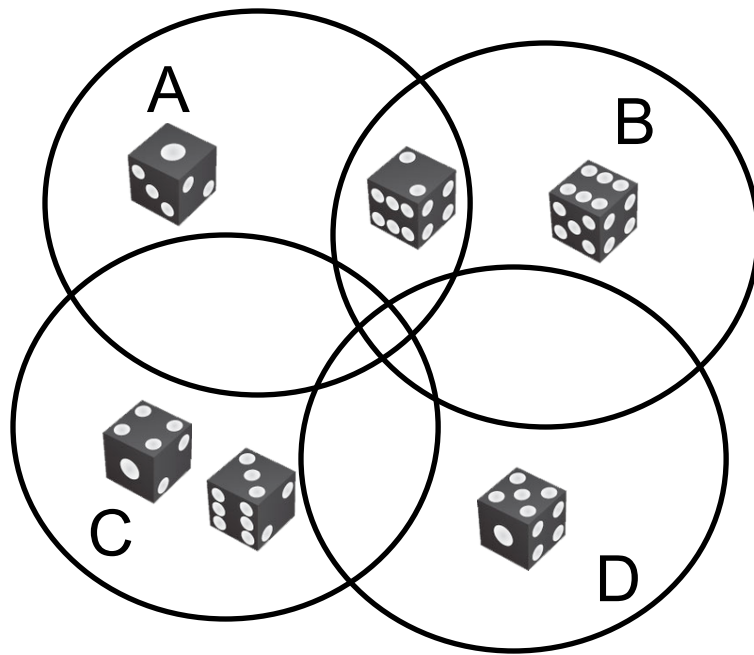
$$O = \{ \text{die with 1, 3, 5} , \text{die with 2, 4, 6} , \text{die with 1, 3, 5} \} \quad \text{Odd die tosses}$$

- The **probability** of an event is given by the sum of the probabilities of the outcomes it contains,

$$p(E) = \sum_{x \in E} p(x) \quad \text{E.g., } p(E) = p(\text{die with 2, 4, 6}) + p(\text{die with 1, 3, 5}) + p(\text{die with 2, 4, 6}) \\ = 1/2, \text{ if fair die}$$

Introduction to probability: union bound

- $P(A \text{ or } B \text{ or } C \text{ or } D \text{ or } \dots)$
 $\leq P(A) + P(B) + P(C) + P(D) + \dots$



$$p(A \cup B) = p(A) + p(B) - p(A \cap B)$$
$$\leq p(A) + p(B)$$

Q: When is this a tight bound?

A: For disjoint events
(i.e., non-overlapping circles)

Introduction to probability: independence

- Two events A and B are **independent** if

$$p(A \cap B) = p(A)p(B)$$

- Suppose our outcome space had two different die:

$$\Omega = \{ \text{die1, die2}, \text{die1, die2}, \text{die1, die2}, \dots, \text{die1, die2} \} \quad \text{2 die tosses}$$

$6^2 = 36$ outcomes

and the probability of each outcome is defined as

$$p(\text{die1, die2}) = a_1 b_1 \quad p(\text{die1, die2}) = a_1 b_2 \quad \dots$$

a_1	a_2	a_3	a_4	a_5	a_6
.1	.12	.18	.2	.1	.3

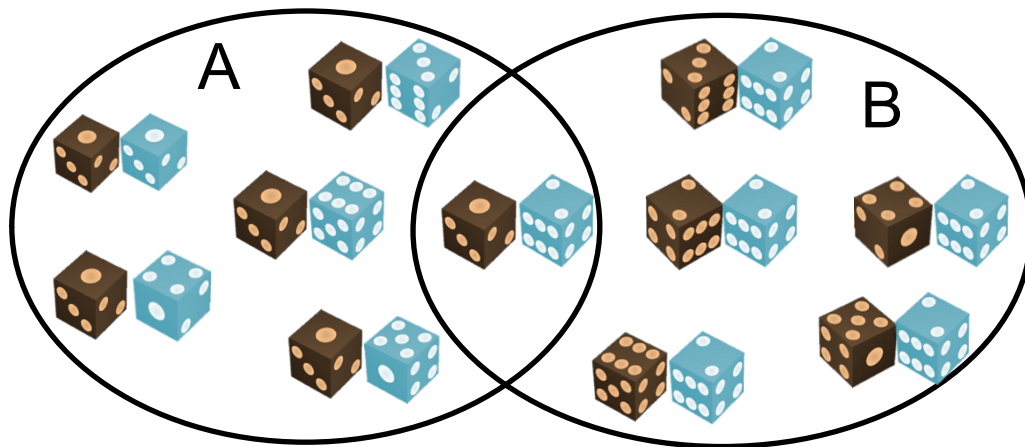
b_1	b_2	b_3	b_4	b_5	b_6
.19	.11	.1	.22	.18	.2

$$\sum_{i=1}^6 a_i = 1$$

$$\sum_{j=1}^6 b_j = 1$$

Introduction to probability: independence

- Two events A and B are **independent** if
$$p(A \cap B) = p(A)p(B)$$
- Are these events independent?



$$p(A) = p(\text{brown die})$$

$$p(B) = p(\text{blue die}) = b_2$$

$$= \sum_{j=1}^6 a_1 b_j = a_1 \sum_{j=1}^6 b_j = a_1$$

Yes! $p(A \cap B) = p(\text{brown die and blue die})$

$$p(A)p(B) = p(\text{brown die}) p(\text{blue die})$$

Introduction to probability: discrete random variables

- A **random variable** X is a mapping $X : \Omega \rightarrow D$
 - D is some set (e.g., the integers)
 - Induces a partition of all outcomes Ω
- For some $x \in D$, we say

$$p(X = x) = p(\{\omega \in \Omega : X(\omega) = x\})$$

“probability that variable X assumes state x ”

- Notation: $\text{Val}(X) = \text{set } D \text{ of all values assumed by } X$
(will interchangeably call these the “values” or “states” of variable X)
- $p(X)$ is a distribution: $\sum_{x \in \text{Val}(X)} p(X = x) = 1$

$$\Omega = \{ \text{🎲🎲}, \text{🎲🎲}, \text{🎲🎲}, \dots, \text{🎲🎲} \} \quad \text{2 die tosses}$$

Introduction to probability: discrete random variables

$X=x$ is simply an event, so can apply union bound, etc.

Two random variables \mathbf{X} and \mathbf{Y} are **independent** if:

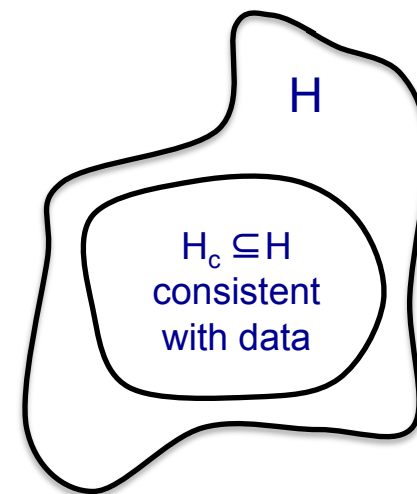
$$p(X = x, Y = y) = p(X = x)p(Y = y) \quad \forall x \in \text{Val}(X), y \in \text{Val}(Y)$$

The **expectation** of \mathbf{X} is defined as: $E[X] = \sum_{x \in \text{Val}(X)} p(X = x)x$

How big should your validation set be?

- In PS1, you tried many configurations of your algorithms (avg vs. regular perceptron, max # of iterations) and chose the one that had smallest validation error
- Suppose in total you tested $|H|=40$ different classifiers on the validation set of m held-out e-mails
- The best classifier obtains 98% accuracy on these m e-mails!!!
- But, what is the true classification accuracy?
- How large does m need to be so that we can guarantee that the best configuration (measured on validate) is truly good?

A simple setting...



- Classification
 - m data points
 - **Finite** number of possible hypothesis (e.g., 40 spam classifiers)
- A learner finds a hypothesis h that is **consistent** with training data
 - Gets zero error in training: $error_{train}(h) = 0$
 - I.e., assume for now that one of the classifiers gets 100% accuracy on the m e-mails (we'll handle the 98% case afterward)
- What is the probability that h has more than ϵ **true** error?
 - $error_{true}(h) \geq \epsilon$

How likely is a **bad** hypothesis to get m data points right?

- Hypothesis h that is **consistent** with validate data
 - got m i.i.d. points right
 - h “bad” if it gets all this data right, but has high true error
 - What is the probability of this happening?
- Probability that h with $\text{error}_{\text{true}}(h) \geq \epsilon$ classifies a randomly drawn data point correctly:
 1. $\Pr(h \text{ gets data point } \textit{wrong} \mid \text{error}_{\text{true}}(h) = \epsilon) = \epsilon$ E.g., probability of a biased coin coming up tails
 2. $\Pr(h \text{ gets data point } \textit{wrong} \mid \text{error}_{\text{true}}(h) \geq \epsilon) \geq \epsilon$
 3. $\Pr(h \text{ gets data point } \textit{right} \mid \text{error}_{\text{true}}(h) \geq \epsilon) = 1 - \Pr(h \text{ gets data point } \textit{wrong} \mid \text{error}_{\text{true}}(h) \geq \epsilon) \leq 1 - \epsilon$
- Probability that h with $\text{error}_{\text{true}}(h) \geq \epsilon$ gets m iid data points correct:
$$\Pr(h \text{ gets } m \textit{ iid} \text{ data points right} \mid \text{error}_{\text{true}}(h) \geq \epsilon) \leq (1 - \epsilon)^m \leq e^{-\epsilon m}$$

E.g., probability of m biased coins coming up heads

Are we done?

$$\Pr(\text{h gets } m \text{ iid data points right} \mid \text{error}_{\text{true}}(\text{h}) \geq \varepsilon) \leq e^{-\varepsilon m}$$

- Says “if h gets m data points correct, then with very high probability (i.e. $1 - e^{-\varepsilon m}$) it is close to perfect (i.e., will have error $\leq \varepsilon$)”
- This only considers **one** hypothesis!
- Suppose 1 billion classifiers were tried, and each was a *random* function
- For m small enough, one of the functions will classify all points correctly – but all have very large true error

How likely is learner to pick a bad hypothesis?

$$\Pr(h \text{ gets } m \text{ iid data points right} \mid \text{error}_{\text{true}}(h) \geq \varepsilon) \leq e^{-\varepsilon m}$$

Suppose there are $|H_c|$ hypotheses consistent with the m data points

- How likely is learner to pick a bad one, i.e. with *true* error $\geq \varepsilon$?
- We need a bound that holds for all of them!

$$P(\text{error}_{\text{true}}(h_1) \geq \varepsilon \text{ OR } \text{error}_{\text{true}}(h_2) \geq \varepsilon \text{ OR } \dots \text{ OR } \text{error}_{\text{true}}(h_{|H_c|}) \geq \varepsilon)$$

$$\leq \sum_k P(\text{error}_{\text{true}}(h_k) \geq \varepsilon)$$

← Union bound

$$\leq \sum_k (1-\varepsilon)^m$$

← bound on individual h_j s

$$\leq |H|(1-\varepsilon)^m$$

← $|H_c| \leq |H|$

$$\leq |H| e^{-m\varepsilon}$$

← $(1-\varepsilon) \leq e^{-\varepsilon}$ for $0 \leq \varepsilon \leq 1$

Generalization error of finite hypothesis spaces [Haussler '88]

We just proved the following result:

Theorem: Hypothesis space H finite, dataset D with m i.i.d. samples, $0 < \epsilon < 1$: for any learned hypothesis h that is consistent on the training data:

$$P(\text{error}_{\text{true}}(h) > \epsilon) \leq |H|e^{-m\epsilon}$$

Using a PAC bound

Typically, 2 use cases:

- 1: Pick ϵ and δ , compute m
- 2: Pick m and δ , compute ϵ

Argument: Since for all h we know that

$$P(\text{error}_{\text{true}}(h) > \epsilon) \leq |H|e^{-m\epsilon}$$

... with probability $1-\delta$ the following holds... (either case 1 or case 2)

$$p(\text{error}_{\text{true}}(h) \geq \epsilon) \leq |H|e^{-m\epsilon} \leq \delta \quad \left. \vphantom{p(\text{error}_{\text{true}}(h) \geq \epsilon)} \right\} \text{ Says: we are willing to tolerate a } \delta \text{ probability of having } \geq \epsilon \text{ error}$$

$\epsilon = \delta = .01, |H| = 40$
Need $m \geq 830$

$$\ln(|H|e^{-m\epsilon}) \leq \ln \delta$$

$$\ln |H| - m\epsilon \leq \ln \delta$$

Case 1

$$m \geq \frac{\ln |H| + \ln \frac{1}{\delta}}{\epsilon}$$

Log dependence on $|H|$, OK if exponential size (but not doubly)

Case 2

$$\epsilon \geq \frac{\ln |H| + \ln \frac{1}{\delta}}{m}$$

ϵ has stronger influence than δ

ϵ shrinks at rate $O(1/m)$

Limitations of Haussler '88 bound

- There may be no consistent hypothesis h (where $error_{train}(h)=0$)
- Size of hypothesis space
 - What if $|H|$ is really big?
 - What if it is continuous?
- **First Goal:** Can we get a bound for a learner with $error_{train}(h)$ in the data set?

Question: What's the expected error of a hypothesis?

- The probability of a hypothesis incorrectly classifying: $\sum_{(\vec{x}, y)} p(\vec{x}, y) 1[h(\vec{x}) \neq y]$
- Let's now let Z_i^h be a random variable that takes two values, 1 if h correctly classifies data point i, and 0 otherwise
- The Z variables are **independent** and **identically distributed** (i.i.d.) with

$$\Pr(Z_i^h = 0) = \sum_{(\vec{x}, y)} p(\vec{x}, y) 1[h(\vec{x}) \neq y]$$

- Estimating the true error probability is like estimating the parameter of a coin!
- Chernoff bound:** for m i.i.d. coin flips, X_1, \dots, X_m , where $X_i \in \{0, 1\}$. For $0 < \epsilon < 1$:

$$p(X_i = 1) = \theta$$

$$P\left(\theta - \frac{1}{m} \sum_i x_i > \epsilon\right) \leq e^{-2m\epsilon^2}$$

True error probability

Observed fraction of points incorrectly classified

$$E\left[\frac{1}{m} \sum_{i=1}^m X_i\right] = \frac{1}{m} \sum_{i=1}^m E[X_i] = \theta$$

(by linearity of expectation)

Generalization bound for $|H|$ hypothesis

Theorem: Hypothesis space H finite, dataset D with m i.i.d. samples, $0 < \epsilon < 1$: for any learned hypothesis h :

$$\Pr(\text{error}_{\text{true}}(h) - \text{error}_D(h) > \epsilon) \leq |H|e^{-2m\epsilon^2}$$

Why? Same reasoning as before. Use the Union bound over individual Chernoff bounds

PAC bound and Bias-Variance tradeoff

for all h , with probability at least $1-\delta$:

$$\text{error}_{\text{true}}(h) \leq \underbrace{\text{error}_D(h)}_{\text{"bias"}} + \underbrace{\sqrt{\frac{\ln |H| + \ln \frac{1}{\delta}}{2m}}}_{\text{"variance"}}$$

- For large $|H|$
 - low bias (assuming we can find a good h)
 - high variance (because bound is looser)
- For small $|H|$
 - high bias (is there a good h ?)
 - low variance (tighter bound)

PAC bound: How much data?

$$\Pr(\text{error}_{\text{true}}(h) - \text{error}_D(h) > \epsilon) \leq |H|e^{-2m\epsilon^2}$$

$$\text{error}_{\text{true}}(h) \leq \text{error}_D(h) + \sqrt{\frac{\ln |H| + \ln \frac{1}{\delta}}{2m}}$$

- Given δ, ϵ how big should m be?

$$m \geq \frac{1}{2\epsilon^2} \left(\ln |H| + \ln \frac{1}{\delta} \right)$$