

STOCHASTIC OPTIMIZATION (AKA MACHINE LEARNING)

15-359 S12

SHIVA KAUL

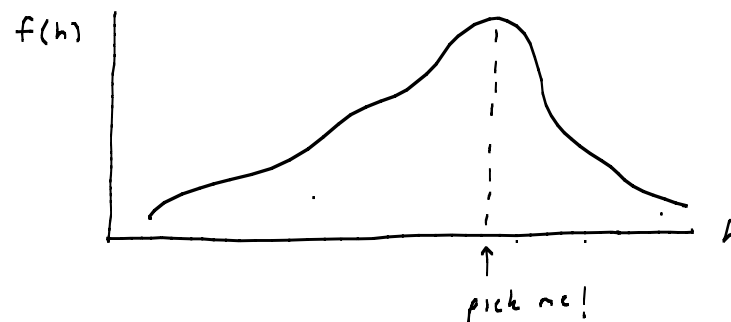
SKKAUL@CS.CMU.EDU

You have probably encountered a deterministic optimization problem...

$$\text{maximize } f(h) \quad \text{s.t. } h \in \mathcal{H}$$

where f is deterministic and \mathcal{H} is some set. This task is computational.

... if not, now you have.



In stochastic optimization:

- f is a RANDOM function drawn from some distribution \mathcal{D} .

The distribution to the right places more probability on the thicker functions.

Now the computational task is a bit more complicated:

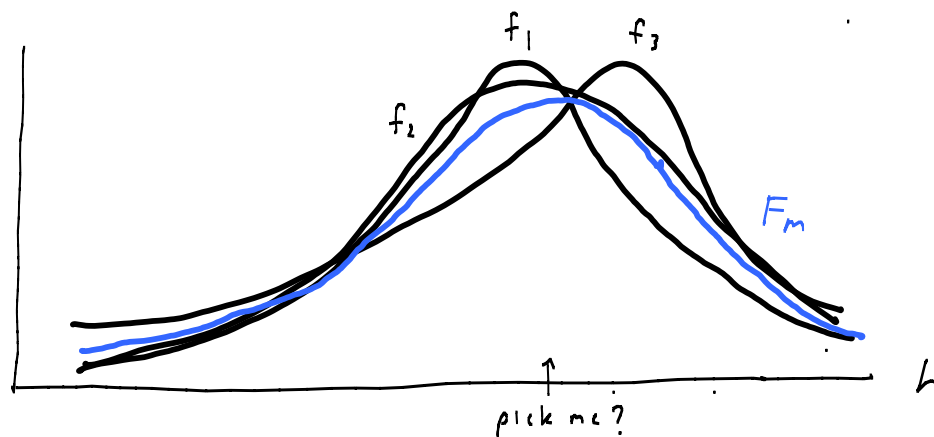
$$\text{maximize } \mathbb{E}_{f \sim \mathcal{D}} [f(h)] \quad \text{s.t. } h \in \mathcal{H}$$



• we don't know the distribution D . All we have is an iid sample

$$f_1, \dots, f_m \sim D$$

now the problem is both computational and statistical ;
we are now estimating the true optimum from limited data.



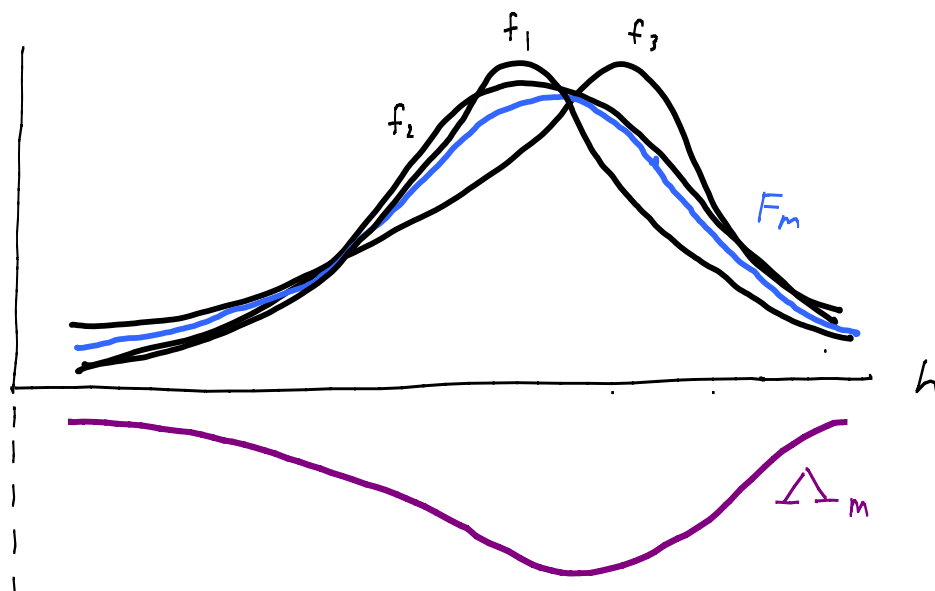
Perhaps a reasonable estimator is the solution to the deterministic problem:

$$\text{maximize } F_m(h) = \frac{1}{m} \sum_{i=1}^m f_i(h) \quad \text{s.t. } h \in \mathcal{H}$$

the 'empirical objective'

We'll study that later. One more change to the optimization setup:

• Due to this uncertainty, we may hedge our estimate by returning a distribution Δ_m over h 's rather than just one h .



this part of
the axis
denotes
probability

Let's overload the notation to accommodate such distributions:

$$f(\Delta_m) = \mathbb{E}_{h \sim \Delta_m} [f(h)]$$

Now let's state the stochastic optimization problem:

given $f_1, \dots, f_m \sim \mathcal{D}$

maximize $F(\Delta_m) = \mathbb{E}_{f \sim \mathcal{D}} [f(\Delta_m)]$ s.t.

In machine learning terminology,

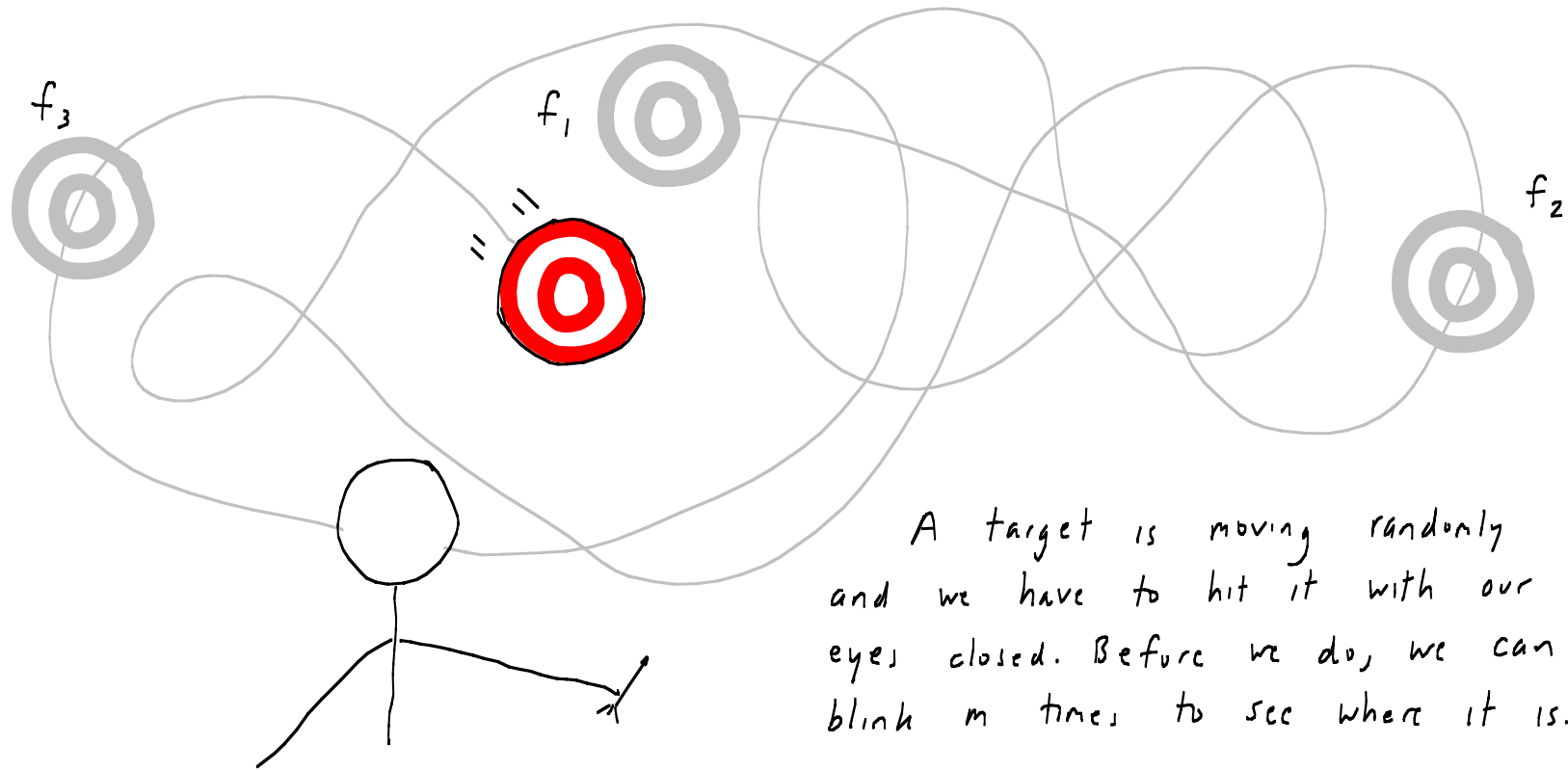
- \mathcal{D} is fixed by nature.
- Each h is a hypothesis about how nature works.

It turns out most of machine learning is stochastic optimization.

But stochastic optimization also includes such important tasks as.....

.... Playing blind carnival games!

Nifty example due
to Geoff Gordon.



A target is moving randomly
and we have to hit it with our
eyes closed. Before we do, we can
blink m times to see where it is.

\mathcal{H} = places to throw dart

$$f^{(x)}(h) = \begin{cases} 1 & \text{if a dart at } h \text{ hits the target at } x \\ 0 & \text{otherwise} \end{cases}$$

\mathcal{D} = distribution over target positions x

----- BEGIN RANT -----

Humans are equipped with many instincts and heuristics which help us model and understand the world. Or, cynically, humans are blinded by vicious biases and prejudices. In pop culture (and unfortunately some pop science), these mental features have corrupted our feeble attempts at rational scientific inquiry, and machine learning is our salvation. Big data, some say, will sweep away our pathetic 'models' and usher in Pax Algorithmica.

Tons of data, fast computers, and clever algorithms are great. Nonetheless, good old-fashioned hunches/intuitions/models are an important and enduring aspect of machine learning.

(Or, cynically, we're still doomed, even with data and computers.)

--- END RANT ---

Our choice of Δ will be good (by the triangle inequality) if:

#1. $F_m(\Delta_m)$ is small (the aforementioned estimator ensures this)

#2. The distance between $F_m(\Delta_m)$ and $F(\Delta_m)$ is small

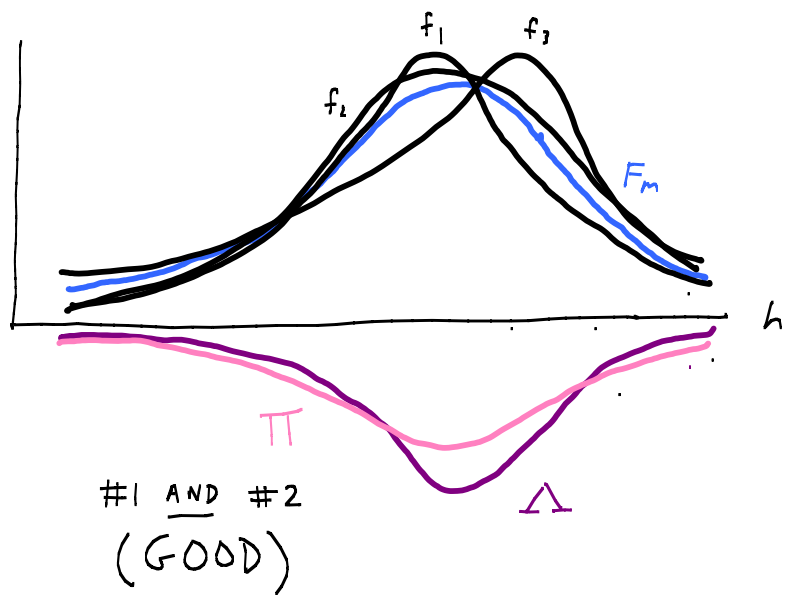
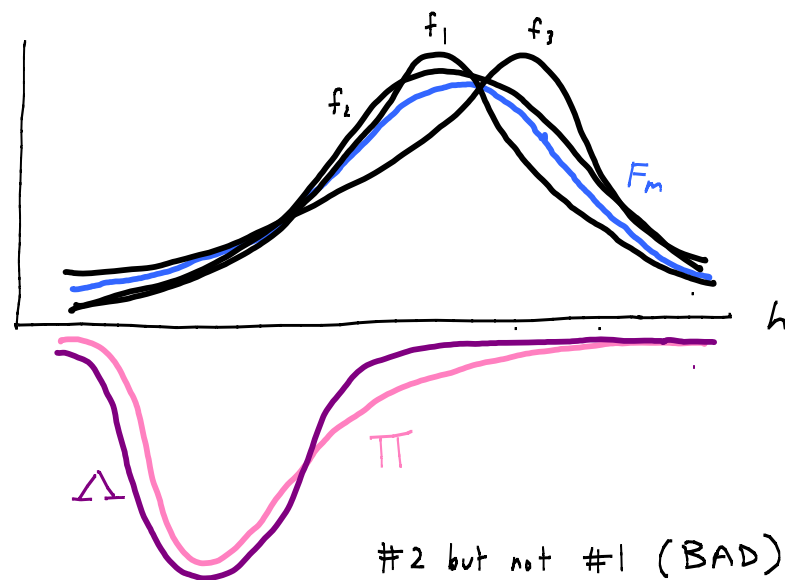
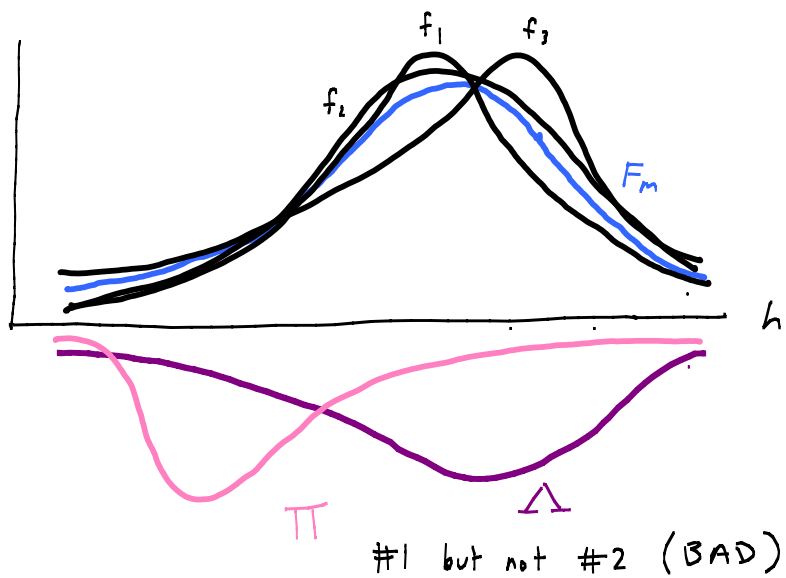
Intuitively, the bigger m is, the better.

The following methodology ^{too abstract to be called an algorithm} highlights the role of hunches.
(We'll describe it conceptually, then quantitatively. Then we'll prove it works)

- start off with a hunch in the form of a distribution Π over H . 'Reasonable' h 's get more probability; 'silly' h 's get less.
- receive the data f_1, \dots, f_m . With probability δ , the following steps don't work.
- for any Δ we choose, these relationships hold (we'll prove this!)

#2 \swarrow distance of $F_m(\Delta_m)$ from $F(\Delta_m)$ \searrow if \cdot distance of Δ_m from Π \downarrow
 \cdot m \uparrow

- strike a balance between #1 and distance of Δ_m from Π .
(a purely computational task)



← a good hunch π is rewarded.

We will quantify the notions of 'distance' or 'closeness' with KL divergence.

$F_m(\Delta)$ vs. $F(\Delta)$:

Remember binary KL-divergence for $q, p \in \{0, 1\}$:

$$kl[q, p] = q \ln\left(\frac{q}{p}\right) + (1-q) \ln\left(\frac{1-q}{1-p}\right)$$

If $\forall f, h, f(h) \in \{0, 1\}$ (like in the dartboard example),
then $F_m(h), F(h) \in \{0, 1\}$, so $kl[F_m(h), F(h)]$ makes sense.

Δ vs. Π :

Remember how $kl[q, p]$ is a special case of a 'distance' (sort of) between distributions? We'll use that.

$$KL[\Delta, \Pi] = E_{h \sim \Delta} \left[\ln\left(\frac{\Delta(h)}{\Pi(h)}\right) \right]$$

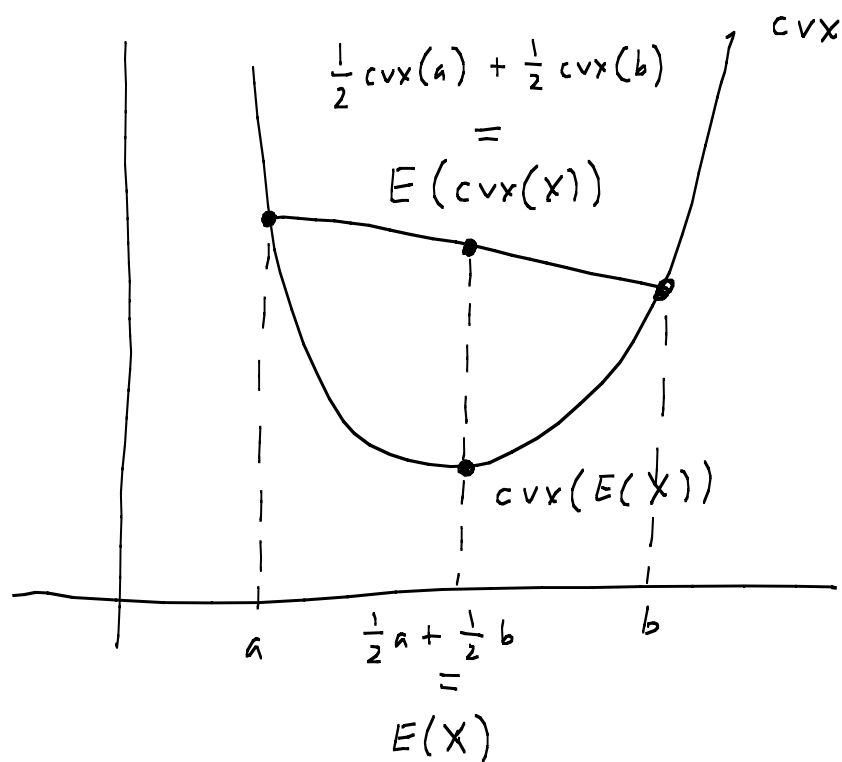
Here's our final methodology, called PAC-Bayes for reasons I won't go into. David McAllester came up with it. I believe this simplified/improved form is by John Langford. ← CMU Grad

- Assumption: $f(h) \in \{0, 1\}$
- start off with a hunch in the form of a distribution π ^{"prior"} over H . 'Reasonable' h 's get more probability; 'silly' h 's get less.
- receive the data f_1, \dots, f_m . With probability δ , the following steps don't work.
- for any Δ ^{"posterior"} we choose,

$$kl \left[F_m(\Delta_m), F(\pi) \right] \leq \frac{KL[\Delta_m, \pi] + \ln \left(\frac{m+1}{\delta} \right)}{m}$$

- strike a balance between $F_m(\Delta_m)$ and $KL[\Delta_m, \pi]$
(a purely computational task)

To prove it, we need Jensen's inequality.



If cvx is a convex function and X is a random variable,

$$cvx(E(X)) \leq E(cvx(X))$$

To remember the direction of the inequality, keep the illustrated special case in mind.



$$X = \begin{cases} a & \text{with probability } 1/2 \\ b & \text{otherwise} \end{cases}$$

Our proof follows the presentation of Ofer Dekel. Define $Z = \mathbb{E}_{h \sim \Pi}(\exp(m \cdot \text{KL}(F_m(h), F(h))))$, a random variable with respect to $S \sim D^m$. We'll prove this theorem in two steps.

$$1. \text{ With probability at least } 1 - \delta, \text{KL}(F_m(\Lambda), F(\Lambda)) \leq \frac{\text{KL}(\Lambda, \Pi) + \ln\left(\frac{\mathbb{E}_{S \sim D^m}(Z)}{\delta}\right)}{m}$$

$$2. \mathbb{E}_{S \sim D^m}(Z) \leq m + 1$$

1. By Markov's inequality,

$$\frac{\mathbb{E}_{S \sim D^m}(Z)}{a} \geq \mathbb{P}(a < Z)$$

Plug in $a = \frac{\mathbb{E}_{S \sim D^m}(Z)}{\delta}$ to get

$$\delta \geq \mathbb{P}\left(\frac{\mathbb{E}_{S \sim D^m}(Z)}{\delta} < Z\right)$$

Inverting this, we have with probability at least $1 - \delta$,

$$\begin{aligned} \frac{\mathbb{E}_{S \sim D^m}(Z)}{\delta} &\geq Z && \text{(above)} \\ \ln\left(\frac{\mathbb{E}_{S \sim D^m}(Z)}{\delta}\right) &\geq \ln(Z) && \text{(monotonicity of log)} \\ &= \ln(\mathbb{E}_{h \sim \Pi}(\exp(m \cdot \text{KL}(F_m(h), F(h)))) && \text{(definition)} \\ &= \ln\left(\mathbb{E}_{h \sim \Lambda}\left(\frac{\Pi(h)}{\Lambda(h)} \exp(m \cdot \text{KL}(F_m(h), F(h)))\right)\right) && \text{(change of measure)} \\ &\geq \mathbb{E}_{h \sim \Lambda}\left(\ln\left(\frac{\Pi(h)}{\Lambda(h)}\right) + m \cdot \text{KL}(F_m(h), F(h))\right) && \text{(concavity of log, Jensen's, log of product)} \\ &= -\text{KL}(\Lambda, \Pi) + m \cdot \mathbb{E}_{h \sim \Lambda}(\text{KL}(F_m(h), F(h))) && \text{(log of quotient, linearity, KL definition)} \\ &\geq -\text{KL}(\Lambda, \Pi) + m \cdot \text{KL}(F_m(\Lambda), F(\Lambda)) && \text{(convexity of KL, Jensen's)} \end{aligned}$$

Rearrange the terms to get part 1. Part 2 is on the homework. \square

Now we have a methodology for stochastic optimization. What do we use it for, and how do we derive algorithms?

An important application is *stochastic supervised learning*, which solves an abstract prediction task: given some observation x (from a ‘feature’ space \mathcal{X}), choose an outcome y (from the outcome space \mathcal{Y}). A loss function $\ell(y, \hat{y}) : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ quantifies how bad it is to choose \hat{y} when the actual outcome is y . (Higher values of ℓ are worse). Some examples:

- Passing/failing P&C: given a student’s homework and quiz grades, predict if they will pass the final exam. Here $\mathcal{X} = \mathbb{R}^d$ where d is the number of assignments/quizzes. $\mathcal{Y} = \{0, 1\}$ encodes ‘fail’ or ‘pass’. $\ell(y, \hat{y}) = 1$ if $y \neq \hat{y}$ and 0 otherwise. We call this setting of \mathcal{Y} and ℓ a classification task.
- Counting faces: given an image, count how many human faces there are. \mathcal{X} is the set of all images, \mathcal{Y} is the set of nonnegative integers, and $\ell(y, \hat{y}) = |y - \hat{y}|$.

In stochastic supervised learning, we assume there is a joint, unknown probability distribution D over the space $\mathcal{X} \times \mathcal{Y}$ (i.e. over random pairs (x, y) .) We seek an $h : \mathcal{X} \rightarrow \mathcal{Y}$ (called a hypothesis) with low risk

$$\ell(h; D) := E_{(x,y) \sim D}(\ell(h; (x, y)))$$

where

$$\ell(h; (x, y)) = \ell(y, h(x))$$

To guide our choice of h , we are given training data:

$$S = \{(x_1, y_1), \dots, (x_m, y_m)\}$$

where each (x_i, y_i) is drawn independently from D . (Note that S is a set, so the order of elements doesn’t matter.) A learning algorithm A takes this training data and returns a hypothesis. This usually involves fitting a hypothesis to the training data - in machine learning parlance, minimizing the empirical risk

$$\ell(h; S) = E_{(x,y) \sim S} \ell(h; (x, y)) = \frac{1}{m} \sum_{i=1}^m \ell(h; (x_i, y_i))$$

As I said, this is a special case of stochastic optimization: take $f(h) = \ell(h; S)$.

How do we derive algorithms? The bound we proved has two especially helpful properties:

- It applies simultaneously to all Λ . A learning algorithm can consider many different Λ and pick one with a good bound.
- If you don't really have a hunch, that's (sometimes) still OK. Choose a Π which isn't too concentrated. When m is small, even if there is an h with small $F_m(h)$, we won't want Λ to put a lot of probability on it, because then the stability term $\text{KL}(\Lambda; \Pi)$ will be large. (Π is shielding us from overfitting.) However, as m increases, the stability term becomes less important, so we can afford to make Λ more concentrated in order to achieve low $F_m(\Lambda)$.