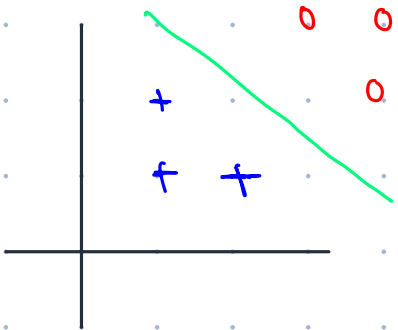# ANNOUNCEMENTS

NOT 11.59 pm

1. HW5 due immediately after spring break (04/07, 5pm)
   (04/09, 5pm — slip days)

2. Tushaar's Th OH cancelled [due to travel], OH over break are cancelled also!

3. Academic Integrity — use of GenAI tools __must__ be disclosed w/ prompts!

4. Other deadlines (P4, P5, etc.), check Ed!

---

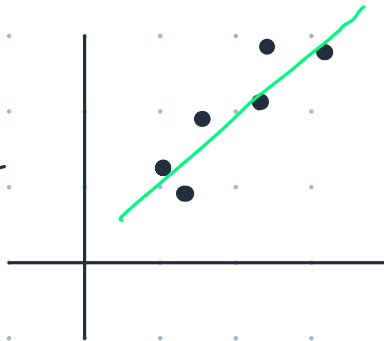As usual, turn your non-note-taking devices off!

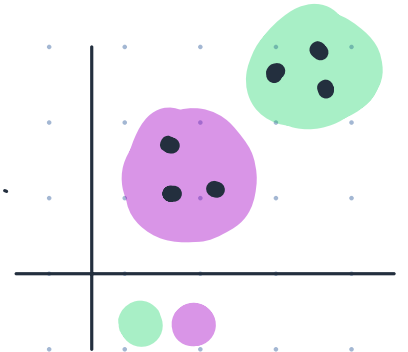So far :      " much time has passed, many algorithms learned "

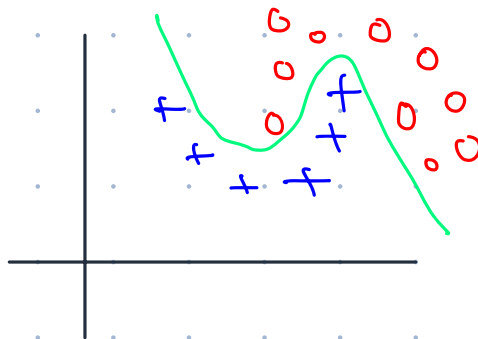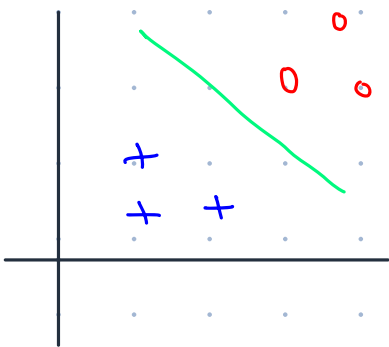( classification )          ( Regression )          Unsupervised

OR

VS.

linearly—separable                  non-linearly- separable

Today — view the fundamental "learning problem"
from the lens of **risk min**

NOT a specific algorithm

Find $\theta$ that minimizes some Cost fn. $J(\theta)$

Q. When does an algorithm / class of algorithm succeed?

# Crickets chirp and temperature raises! — <mark>"generalization"</mark>

temperature

Chirp rate

$temp'' \propto$ chirp rate

temperature

Chirp rate

$temp'' = \theta_0 + \theta_1 \; rate + \theta_2 \; rate^2$

temperature

Chirp rate

griffin says too specific to data

# Same game, different name — Classification w/ logistic regression

# A case of binary classification

Given $D = \{(x^{(j)}, y^{(j)})\}_{j=1}^{n}$ of $n$ samples.
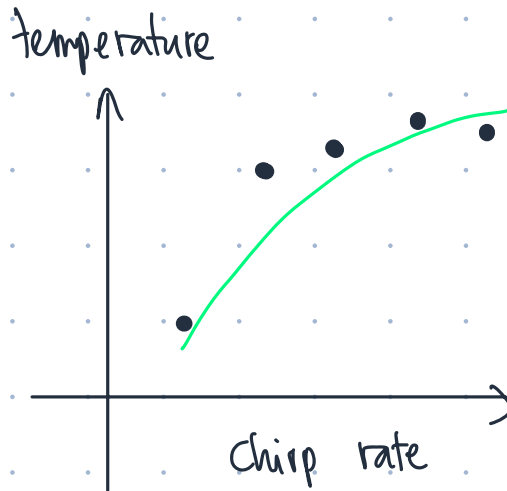
$$(x^{(j)}, y^{(j)}) \overset{iid}{\sim} P$$

Given a fn, hypothesis "$h$", the training error

$$\hat{\varepsilon}(h) = \frac{1}{n} \sum_{j=1}^{n} \boxed{1\{h(x^{(j)}) \neq y^{(j)}\}} \longrightarrow \text{misclassification}$$

train error OR
Empirical risk

$$\hat{\varepsilon}_D(h)$$

$\downarrow$ prediction

We want to choose "$\hat{\theta}$" s.t.

$$\hat{\theta} = \arg\min_{\theta} \hat{\varepsilon}(h)$$

$\Big]$ — Empirical risk minimization!

Logistic regression, SVMs, etc. are Convex approximates
to 0/1 loss

$\max\left(0, 1 - \theta^T x\right)$

$\log\left(1 + \exp(-\theta^T x)\right)$

"loss" or "error" of misclassifying
$y = +1$ points

$\theta^T x$

From "find the optimal $\theta$ to minimize $\hat{\mathcal{E}}(h_\theta)$" to
   Empirical risk minimization

given   $\mathcal{H}$, a class of all hypotheses,

$$\hat{h} = \arg\min\limits_{h \in \mathcal{H}} \hat{\mathcal{E}}(h)$$

linear : $\mathcal{H} = \left\{ x \to \sigma(\theta^T x) \mid \theta \in \mathbb{R}^{d+1} \right\}$    |    linear : $\mathcal{H} = \left\{ x \to \theta^T x \mid \theta \in \mathbb{R}^{d+1} \right\}$
classifiers                                                                                   regression


What do we care about?        NOT  $\hat{\mathcal{E}}(h)$ = train error

   Generalization error

$$\mathcal{E}(h) = P_{(x,y) \sim P} \; \mathbb{1}\left\{ h(x) \neq y \right\}$$

   probability of misclassifying some "$x$"

# What do we want?

1) given the training error, $\hat{\mathcal{E}}(h)$ can we guarantee anything about generalization? $\mathcal{E}(h)$

2) Can we estimate this "generalization error" — $\mathcal{E}(h)$

if the traing dataset is ___ at least large,

It at most ___ complex,

then ___ w.p. , we have training error within $\tau$ of test error.

# Some preliminaries

Lemma

(The union bound).  $\overset{\ell^{\,"or"}}{P(A_1 \cup A_2 \ldots \cup A_k)} \leq P(A_1) + P(A_2) + \ldots + P(A_k)$

$A_1, \ldots, A_k$  events, not necessarily independent



Lemma

(The Chernoff bound).   If we had a coin w/ $P(A) = p$, flipped "$n$" times, then $\hat{p} =$ fraction of times we see $A$

$$P\left( |p - \hat{p}| > r \right) \leq 2 \exp(-2r^2 n)$$

$\hat{p}$

$\leq 2 \exp(-2r^2 n)$

Takeaway: As $n \uparrow$, we get **better** estimates for $\hat{p}$

"$\exp$" better



0          $p$      1

# Relate train error to generalization error

$1\{h(x) \neq y\} = z$ — $(x, y)$ misclassified?       $\mathcal{E}(h) =$ generalization

$\quad P(z = 1) = \mathcal{E}(h) \longrightarrow p$ in coin toss       $\hat{\mathcal{E}}(h) =$ train

$\qquad\qquad\qquad\qquad\qquad\qquad$ example

$1\{h(x^{(j)}) \neq y^{(j)}\} = z^{(j)}$ — $(x^{(j)}, y^{(j)})$ misclassified

$\quad$ now, $\quad \dfrac{1}{n} \displaystyle\sum_{j=1}^{n} z^{(j)} = \hat{\mathcal{E}}(h) \longrightarrow \hat{p}$ in coin toss

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ example

$$P\left( \left| \mathcal{E}(h) - \hat{\mathcal{E}}(h) \right| > \gamma \right) \leq 2 \exp\left( -2\gamma^2 n \right)$$

$\qquad\qquad$ for $\underline{ONE}$ $h \in \mathcal{H}$

$$P\left( \ |\mathcal{E}(h) - \hat{\mathcal{E}}(h)| > r \right) \leq 2\exp\left(-2r^2 n\right)$$

for **ONE** $h \in \mathcal{H}$

Wish to extend to any $h \in \mathcal{H}$ :

Q. what is the probability that one or more $h_j$s result in $|\mathcal{E}(h) - \hat{\mathcal{E}}(h)| > r$ ?

Let $A_j$ to be the event that $\left( \mathcal{E}(h_j) - \hat{\mathcal{E}}(h_j) \right| > r$

$$P\left( A_1 \cup A_2 \cup \ldots \cup A_k \right) \leq P(A_1) + \ldots + P(A_k)$$

$$\leq \sum_{j=1}^{k} 2\exp\left(-2r^2 n\right)$$

$$= 2k \exp\left(-2r^2 n\right)$$

← error margin that we chose

$|\mathcal{H}|$ = num of hypothesis     n = dataset size

$$P\left( \neg\, h_j \in \mathcal{H} \ \text{s.t.} \ |\mathcal{E}(h_j) - \hat{\mathcal{E}}(h_j)| > r \right) \geq 1 - 2k\exp\left(-2r^2 n\right)$$

probability that **none** of $h_1, \ldots, h_k$ result in $|\mathcal{E}(h_j) - \hat{\mathcal{E}}(h_j)| > r$

Alternatively, …

Given some $r > 0$, $0 < \delta < 1$, how large a dataset is needed to guarantee w.p. $1-\delta$, the training error is within "$r$" of generalization error?

Previous result —

$$P\left(\neg\, h_j \in \mathcal{H} \text{ s.t. } |\varepsilon(h_j) - \hat{\varepsilon}(h_j)| > r\right) \geq 1 - 2k \exp\left(-2r^2 n\right)$$

probability that __none__ of $h_1, \dots, h_k$
result in $|\varepsilon(h_j) - \hat{\varepsilon}(h_j)| > r$

$$1 - \delta \geq 1 - 2k \exp\left(-2r^2 n\right)$$

$$\delta \leq 2k \exp\left(-2r^2 n\right)$$

$$\log \frac{\delta}{2k} \leq -2r^2 n \qquad \Rightarrow -2r^2 n \geq \log \frac{\delta}{2k}$$

$$n \geq -\frac{1}{2r^2} \log \frac{\delta}{2k}$$

$$(\text{or}) \quad n \geq \frac{1}{2r^2} \log \frac{2k}{\delta}$$

$$n = O_{r,\delta}\left(\log k\right)$$

## Error bound

Fix $n$ and $\delta$, solve for margin —

$$\left| \mathcal{E}(h_j) - \hat{\mathcal{E}}(h_j) \right| \leq \sqrt{\frac{1}{2n} \log \frac{2k}{\delta}}$$

So, what can we say about $\hat{h}$?

ERM — identifies $\hat{h} = \underset{h \in \mathcal{H}}{\arg\min} \, \hat{\mathcal{E}}(h)$ → best ERM hypothesis

$h^* = \underset{h \in \mathcal{H}}{\arg\min} \, \mathcal{E}(h)$ → best-in-$\mathcal{H}$ hypothesis

$\mathcal{E}(\hat{h}) \le \hat{\mathcal{E}}(\hat{h}) + r$

$\le \hat{\mathcal{E}}(h^*) + r$ → Since $\hat{h}$ is chosen to minimize train err $\hat{\mathcal{E}}$,

$r = 2\sqrt{\ldots}$

no $h^*$ can achieve lower tr error

$= \boxed{\hat{\mathcal{E}}(h^*) - \mathcal{E}(h^*)} + \mathcal{E}(h^*) + r$

$|\hat{\mathcal{E}}(h^*) - \mathcal{E}(h^*)| \le r$

$\le r + \mathcal{E}(h^*) + r$

OR $\quad$ $\mathcal{E}(\hat{h}) \le \mathcal{E}(h^*) + 2r$

ORACLE INEQUALITY!

# Finite $\mathcal{H}$?

1) Deepseek $-671B$ model, $2^{32 \times 671B}$ hypotheses

We need $190T$ examples, to say with $50\%$ proba that generalization error of my model is $20\%$ worse that the best model

# Model selection

    — Can we estimate $\mathcal{E}(h)$ somehow?
         what if $\hat{\mathcal{E}}(h)$ chosen to be $\mathcal{E}(h)$?

① 70/30 split, estimate $\mathcal{E}(h)$ using 30% split

② K-fold LV —



③ leave-one-out    — every data point is used to test!