# Bayes Decision Theory

How to make decisions in the presence of uncertainty?

History: $2^{nd}$ World War

Radar for detection aircraft.

Codebreaking. Decryption.

Observed Data $x \in X$

State $y \in Y$.    <u>likelihood function</u>

$p(x|y)$ — conditional distribution
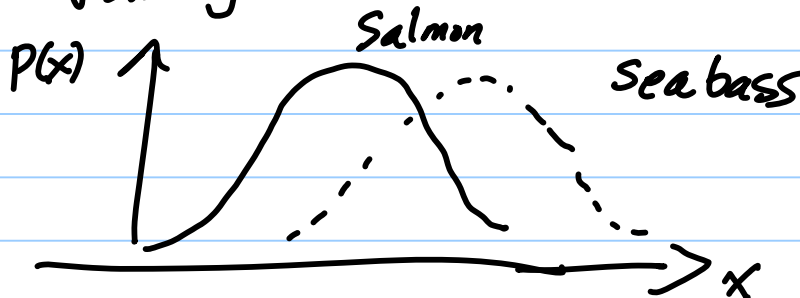model how data is generated.

Example $y \in \{-1, 1\}$    Salmon / Sea Bass
Airplane / Bird

$$p(x|y) = \frac{1}{\sqrt{2\pi}\, \sigma_y} e^{-\frac{1}{2}\frac{(x-\mu_y)^2}{2\sigma_y^2}}$$

mean $\mu_y$:
variance $\sigma_y^2$.

Eg. $x$ is length of fish.

$p(x)$    Salmon    Sea bass

(2)   How to decide   Sea Bass or Salmon?  <span style="color:red">Spring 2013</span>

Airplane or Bird

Maximum Likelihood (ML)

$$\hat{y}_{ML} = \underset{y}{ARG\ MAX}\ p(x|y)$$

$$\left( p(x|\hat{y}_{ML}) \geq p(x|y) \quad \forall y \right)$$

If $\quad P(x|y=1) > P(x|y=-1)$   decide $y=1$

otherwise $y=-1$

Equivalently $\quad \log \dfrac{P(x|y=1)}{P(x|y=-1)} > 0$   log-likelihood test.

Seems reasonable, but what if birds are more likely than airplanes?

Must take into account the _prior probability_ $\quad P(y=1),\ P(y=-1)$.

Bayes Rule $\quad P(y|x) = \dfrac{P(x|y)P(y)}{P(x)}$

Prob of $y$ conditioned on observation.

If $\quad p(y=1|x) > p(y=-1|x)$ decide $y=1$

otherwise decide $y=-1$

Maximum a Posteriori (MAP) $\quad \hat{y}_{MAP} = \underset{y}{ARG\ MAX}\ p(y|x)$

(3)    <u>Another ingredient</u>

→ what does it cost if you make a mistake?

i.e. suppose you decide $y=1$, but really $y=-1$.

i.e. you may pay a big penalty if you decide it is a bird when it is a plane.

(Pascal's Wager: Bet on God)

Putting everything together.

likelihood function   $p(x|y)$      $x \in X, y \in Y$

prior          $p(y)$

decision rule   $\alpha(x)$              $\alpha(x) \in Y$

loss function   $L(\alpha(x), y)$    Cost of making decision $\alpha(x)$ when true state is $Y$.

<u>E.G.</u>   $L(\alpha(x), y) = 0$, if $\alpha(x) = y$
       $L(\alpha(x), y) = 1$, if $\alpha(x) \neq y$

All wrong answers penalized the same.

# (4)   Risk

The <u>risk</u> of the decision rule $d(x)$ is the <u>expected loss</u>.

$$R(d) = \sum_{x,y} L(d(x),y) \, P(x,y)$$

(Note integrate $\int dx$ if $x$ is continuous)

Bayes Decision Theory says "pick the decision rule $\hat{d}$ which minimizes the risk".

$$\hat{d} = \underset{d \in A}{\text{ARGMIN}} \; R(d), \quad R(\hat{d}) \geqslant R(d) \quad \forall \, d \in A.$$

A = set of all decision rules

$\hat{d}$ is Bayes Decision
$R(\hat{d})$ is Bayes Risk.

(5)  # Bayes Risk

Bayes Risk is the best you can do if: (a) you know $p(x|y) p(y)$ & $L(\cdot, \cdot)$

(b) you can compute $\hat{\alpha} = \text{ARGMIN}_{\alpha} R(\alpha)$

(c) you can afford the losses (e.g. gambling, poker)

(d) you make the decision for a sequence of data $x_1 \dots x_n$ with states $y_1 \dots y_n$ where each $(x_i, y_i)$ are independently identically distributed from $p(x,y)$

---

Bad — if you are playing a game against an intelligent opponent (Game Theory)

— if any of the assumptions (a), (b), (c) (d) are wrong.

Note: Cognitive Scientists have studied decision theory to see if it predicts the way human's make decisions. Results are debatle. But Prospect Theory (Kahneman, Trevsky) suggests that humans do not.

6) <u>Better</u> understanding of <span></span>

Bayes Decision Theory.  Re-express

$$R(\alpha) = \sum_x \sum_y L(\alpha(x), y) P(x, y)$$

$$= \sum_x P(x) \left\{ \sum_y L(\alpha(x), y) P(y|x) \right\}$$

Hence, for each $x$,

$$\hat{\alpha}(x) = \underset{\alpha(x)}{\text{ARG MIN}} \sum_y L(\alpha(x), y) P(y|x)$$

<u>Obtaining MAP & ML as special cases</u>.

If $y \in \{-1, 1\}$ and the loss function penalizes all errors equally:

$$L(\alpha(x), y) = 1, \text{ if } \alpha(x) \neq y.$$
$$= 0, \text{ otherwise}$$

$$y \in \{-1, 1\}$$

Then $\hat{\alpha}(x) = \underset{\alpha(x)}{\text{ARG MAX}} P(y = \alpha(x)|x)$

MAP estimate.

If also $p(y=1) = p(y=-1)$, then

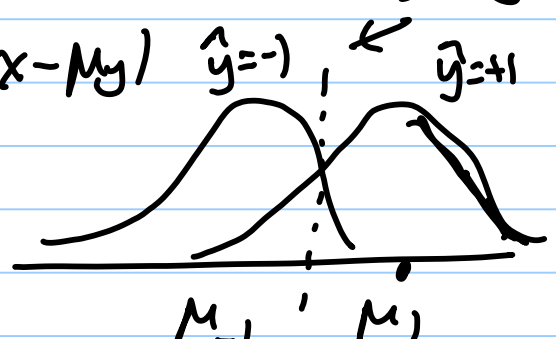$$\hat{\alpha}(x) = \underset{\alpha(x)}{\text{ARG MAX}} p(x|y=\alpha(x)) \text{ ML estimate}$$

(7)  Examples
$$p(x|y) = \frac{1}{\sqrt{2\pi}\,\sigma_y} e^{-\frac{(x-\mu_y)^2}{2\sigma^2}}$$

$y \in \langle -1, 1 \rangle \qquad p(y) = \frac{1}{2}$

$L(\alpha(x), y) = 1, \text{ if } \alpha(x) \neq y, \quad = 0 \text{ otherwise.}$

Bayes Rule
$$\alpha(x) = \underset{y \in \{-1, 1\}}{\text{ARG MIN}} |x - \mu_y|$$



Decision Boundary

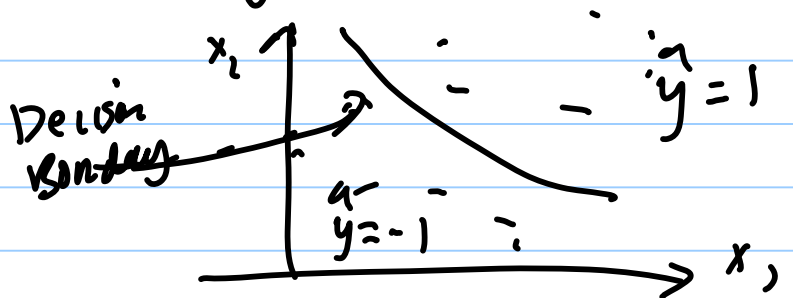$\hat{y} = -1 \qquad \hat{y} = +1$

$\mu_{-1} \qquad \mu_1$

Suppose $\underline{x}$ is a vector in two dimensions

$$p(\underline{x}|y) = \frac{1}{2\pi\sigma^2} e^{-\frac{1}{2\sigma^2}|\underline{x}-\underline{\mu}_y|^2}$$

Separating Plane ↔ Decision Boundary



$x_2$

$\hat{y} = -1$

$\hat{y} = +1$

$x_1$

$$\text{If } p(\underline{x}|y) = \frac{1}{2\pi|\Sigma_y|^{\frac{1}{2}}} e^{-\frac{1}{2}(\underline{x}-\underline{\mu}_y)^T \Sigma_y^{-1} (\underline{x}-\underline{\mu}_y)}$$

Gaussians with unequal covariances

Decision Boundary



$x_2$

$\hat{y} = 1$

$\hat{y} = -1$

$x_1$

<u>More Details</u>

$$P(\underline{x}|y) = \frac{1}{2\pi|\Sigma|^{\frac{1}{2}}} e^{-\frac{1}{2}(\underline{x}-\underline{\mu}_y)^T \Sigma^{-1}(\underline{x}-\underline{\mu}_y)}$$

ie same covariance $\Sigma$ for both classes $y = \pm1$.

$$\log \frac{P(\underline{x}|y=1)}{P(\underline{x}|y=-1)} = \frac{1}{2}(\underline{x}-\underline{\mu}_{-1})^T \Sigma^{-1}(\underline{x}-\underline{\mu}_{-1})$$
$$-\frac{1}{2}(\underline{x}-\underline{\mu}_1)^T \Sigma^{-1}(\underline{x}-\underline{\mu}_1) \qquad \left(\begin{array}{c}2\pi|\Sigma|^{\frac{1}{2}} \text{ terms} \\ \text{cancel}\end{array}\right)$$

$$= (\underline{\mu}_{-1}-\underline{\mu}_1)^T \Sigma^{-1}\underline{x} + \frac{1}{2}\underline{\mu}_{-1}^T \Sigma^{-1}\underline{\mu}_{-1}$$
$$-\frac{1}{2}\underline{\mu}_1^T \Sigma^{-1}\underline{\mu}_1.$$

Linear in $\underline{x}$ describes a <u>plane</u>.

Hence ML rule/estimator corresponds to a rule.

Classify $\underline{x}$ as $y=1$ if
$$(\underline{\mu}_{-1}-\underline{\mu}_1)^T \Sigma^{-1}\underline{x} + \frac{1}{2}\underline{\mu}_{-1}^T \Sigma^{-1}\underline{\mu}_{-1} - \frac{1}{2}\underline{\mu}_1^T \Sigma^{-1}\underline{\mu}_1 > 0$$

as $y=-1$ if $\quad (\underline{\mu}_{-1}-\underline{\mu}_1)^T \Sigma^{-1}\underline{x} + \frac{1}{2}\underline{\mu}_{-1}^T \Sigma^{-1}\underline{\mu}_{-1} - \frac{1}{2}\underline{\mu}_1^T \Sigma^{-1}\underline{\mu}_1 < 0.$

If there is a prior $p(y)$

$$\log \frac{P(y=1|\underline{x})}{P(y=-1|\underline{x})} = \log \frac{P(\underline{x}|y=1) P(y=1)}{P(\underline{x}|y=-1) P(y=-1)}$$

$$= \log \frac{P(\underline{x}|y=1)}{P(\underline{x}|y=-1)} + \log \frac{P(y=1)}{P(y=-1)} \quad \longleftarrow \text{ Indep of } \underline{x}$$

$\left(\begin{array}{c}P(y|\underline{x}) = \frac{P(\underline{x}|y)P(y)}{P(\underline{x})} \\ P(\underline{x}) \text{ cancels} \\ \text{in the ratio.}\end{array}\right)$

<u>Hence</u> prior shifts the seperating plane to
$$(\underline{\mu}_{-1}-\underline{\mu}_1)^T \Sigma^{-1}\underline{x} + \frac{1}{2}\underline{\mu}_{-1}^T \Sigma^{-1}\underline{\mu}_{-1} - \frac{1}{2}\underline{\mu}_1^T \Sigma^{-1}\underline{\mu}_1 + \log\frac{P(y=1)}{P(y=-1)}.$$

<u>With Loss Function</u> $R(\alpha(\underline{x})=1) = L(1,1)P(y=1|\underline{x}) + L(1,-1)P(y=-1|\underline{x})$
$R(\alpha(\underline{x})=-1) = L(-1,1)P(y=1|\underline{x}) + L(-1,-1)P(y=-1|\underline{x})$

Decision boundary occurs where $R(\alpha(\underline{x})=1) = R(\alpha(\underline{x})=-1).$
ie, when $\langle L(1,1) - L(-1,1)\rangle P(y=1|\underline{x}) = \langle L(-1,-1) - L(1,-1)\rangle P(y=-1|\underline{x})$
ie when $\log \frac{P(y=1|\underline{x})}{P(y=-1|\underline{x})} = \log \frac{\langle L(-1,-1) - L(1,-1)\rangle}{\langle L(1,1) - L(-1,1)\rangle}$

$\nearrow$ additional shift in position of seperating plane.

Bayes Decision theory also applies when $y$ is not a binary variable — e.g. $y$ can take M values or $y$ continuous valued.

In this Course, usually $y \in \langle -1, 1 \rangle$ classification

or $y \in \langle 1, 2 \ldots M \rangle$ multi-class classification

or $y \in \{ -\infty, \infty \rangle$ regression.

But there is a major problem with Bayes Decision Theory, apart from the limitations discussed earlier.

Problem, We almost never know the probability distributions $p(y|x)$ and $p(\underline{x})$.

Instead we have data $\langle (\underline{x}_i, y_i) : i = 1 \text{ to } N \rangle$

E.G. Bank has records of the incomes and savings of the customers — $\underline{x}_i$. and whether they defaulted on their loans — $y_i$.

Similarly, you can make a dataset of fish, recording their length and brightness $\underline{x}_i$ and whether they are salmon or sea bass $y_i$

## Two Strategies

## (1) Probability Approach.

Use the data $\langle (x_i, y_i) : i = 1 \text{ to } N \rangle$ to learn probability
distribution $P(x|y)$ and $p(y)$. Then apply Bayes Decision Theory.

E.G. $\quad p(y=1) = \sum_{i=1}^{N} \dfrac{I(y_i=1)}{N}$ $\qquad$ Indicator function:

$\qquad\qquad\qquad P(y=-1) = \sum_{i=1}^{N} \dfrac{I(y_i=-1)}{N}$ $\qquad$ $I(y=1) = 1$, if $y=1$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad I(y=1) = 0$, otherwise

Gaussian assumption $\quad P(x|y=1) = N(\mu_1, \Sigma_1)$ $\qquad$ $N(\cdot)$ mean

$\qquad\qquad\qquad\qquad\qquad P(x|y=-1) = N(\mu_{-1}, \Sigma_{-1})$ $\qquad$ normal / covariance

$$\mu_1 = \frac{\sum_{i=1}^{N} I(y_i=1) x_i}{\sum_{i=1}^{N} I(y_i=1)} \quad , \quad \mu_{-1} = \frac{\sum_{i=1}^{N} I(y_i=-1) x_i}{\sum_{i=1}^{N} I(y_i=-1)}$$

$$\Sigma_1 = \frac{1}{\sum_{i=1}^{N} I(y_i=1)} \sum_{i=1}^{N} I(y_i=1)(x_i - \mu_1)(x_i - \mu_1)^T$$

$$\Sigma_2 = \frac{1}{\sum_{i=1}^{N} I(y_i=-1)} \sum_{i=1}^{N} I(y_i=-1)(x_i - \mu_{-1})(x_i - \mu_{-1})^T$$

I.e. estimate the mean and covariances for
classes $y=1$ and $y=-1$ using only the data assigned
to that class (e.g. assign $x_i$ to class $y=1$, if $y_i=1$).

Note: This strategy requires learning
parametric and non-parametric probability distribution
— we will discuss methods for doing this
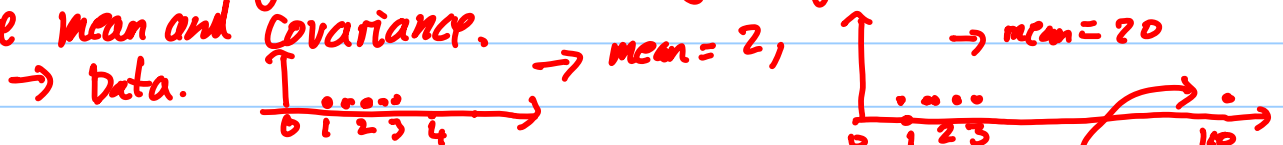in later lectures.

<u>Strategy (2)</u>

<u>Discriminative</u> : Attempt to learn the decision rule
$d(x)$ directly from the data $\langle (\underline{x}_i, y_i) : i \in 1 6 N \rangle$.

i.e. select $\hat{d}(.) = $ ARG MIN $R_{emp}(d ; N)$
$$d(.) \quad = \text{ARG MIN} \sum_{i=1}^{N} \mathcal{L}(d(\underline{x}_i), y_i)$$
$$d(.)$$

<u>Justification</u> — why bother learning the probabilities if
we really only care about the decision rule.?

<u>E.G.</u> Suppose we use Gaussians to model $P(\underline{x}|y)$.
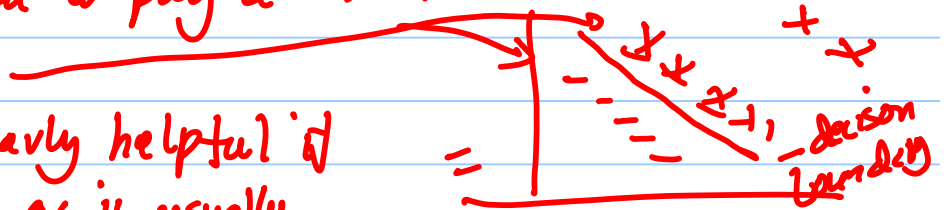
Gaussians are <u>non-robust</u> — outliers in the data, i.e.
atypical values of $\underline{x}_i$, can make big changes to the estimates
of the mean and covariance.
→ Data.

→ mean = 2,  → mean = 20

one outlier

So if we to learn $P(x|y)$ our estimates
of the mean and covariance, hence of the decision boundary,
can be corrupted by outliers far away from the boundary.

But if instead, we just search for a linear
plane that seperates the data from $y=1$ and $y=-1$
then we only need to pay attention to the data near
the boundary

This is particularly helpful if
we have little data, as is usually
the case.

## Key Issue of Machine Learning:

How to generalize — learn a classification rule on training data $\{(\underline{x}_i, y_i) : i = 1 \text{ to } N\}$ which also works on new data that you haven't seen. i.e. predict on new data.

How to formalize this?

Assume that the training data $\{(\underline{x}_i, y_i) : i = 1 \text{ to } N\}$ are independently identically distributed (i.i.d) from an unknown distribution $P(\underline{x}, y)$.

Want the decision rule $\alpha(.)$ trained on $\{(\underline{x}_i, y_i) : i = 1 \text{ to } N\}$ to work on other i.i.d samples $\{(\underline{x}_j, y_j) : j = N+1 \text{ to } N+M\}$, from $P(\underline{x}, y)$

i.e. if empirical risk $R_{emp}(\alpha : N)$ is small then want the risk $R(\alpha)$ to be small.

Now $\alpha \in A$, where $A$ is a set of decision rules (eg. could include ML, MAP separating planes, nearest neighbor, decision trees).

So if we can make sure that
$$|R_{emp}(\alpha, N) - R(\alpha)| \text{ is small for all } \alpha \in A$$
then we can select a rule $\tilde{\alpha} = \arg\min_{\alpha} R(\alpha : N)$ and be confident that
$$R(\tilde{\alpha}) \text{ is close to } \min_{\alpha} R(\alpha)$$
ie that rule $\tilde{\alpha}$ generalizes.

Memorization vis. Generalization

$$R_{emp}(\alpha) = \frac{1}{N} \sum_{i=1}^{N} L(\alpha(x_i), y_i)$$

$$\text{suppose } L(\alpha(x_i), y_i) \in \left\{ \begin{matrix} 1 \\ 0 \end{matrix} \right\}$$

By law of large numbers $R_{emp}(\alpha) \xrightarrow[N \to \infty]{} R(\alpha) = \sum_{x,y} p(x,y) L(\alpha(x),y)$
but how fast?

Fix $\alpha$: By standard theorems (Chernoff, Sanov, Cramers.)

$$Pr\left\{ |R_{emp}(\alpha) - R(\alpha)| > \epsilon \right\} < e^{-N\epsilon}$$

require $e^{-N\epsilon} < \delta \iff N > -\frac{1}{\epsilon} \log \delta$  $\left\{ \begin{matrix} \text{any } \epsilon \end{matrix} \right\}$  $\left( \begin{matrix} -\log \delta > 0 \\ \text{if } 0 \leq \delta \leq 1 \end{matrix} \right)$

So, if $N > -\frac{1}{\epsilon} \log \delta$, then with prob $> 1 - \delta$

$$|R_{emp}(\alpha) - R(\alpha)| < \epsilon.$$  Almost sure that we can.

**Probably Approximately Correct (PAC)** estimate Bayes risk from $N$ samples.

But, we have to consider many different rules $\alpha$. For simplicity, suppose we consider a finite no. of rules $\left\{ \alpha^\nu : \nu = 1 \text{ to } H \right\}$.

We want $|R_{emp}(\alpha^\nu) - R(\alpha^\nu)| < \epsilon$ to be small for all $\nu$. with high probability.

Boole's Inequality: $Pr(A^1 \text{ or } \dots \text{ or } A^H) \leq \sum_{\nu=1}^{H} Pr(A^\nu)$

Let $Pr(A^\nu)$ be prob that $|R_{emp}(\alpha^\nu) - R(\alpha^\nu)| > \epsilon$

$Pr\left\{ \text{At least one rule } A^\nu \text{ has error greater than } \epsilon \right\}$

$$< H e^{-N\epsilon}$$  Now want $H e^{-N\epsilon} < \delta$

$$\iff N > \frac{1}{\epsilon} \left\{ \log H - \log \delta \right\}$$

So if $N > \frac{1}{\epsilon} \left\{ \log H - \log \delta \right\}$, then with prob $> 1 - \delta$

$$|R_{emp}(\alpha^\nu) - R(\alpha^\nu)| < \epsilon \text{ for all } \nu = 1 \text{ to } H.$$  size of hypothesis space

Hence number of examples needed grows rapidly with $H$, accuracy required $\epsilon$, certainty $\delta$.

## Memorization:

Decision Rule: $\hat{\alpha} = \underset{\alpha}{ARGMIN}\ R_{emp}(\alpha)$

$R_{emp}(\hat{\alpha})$ small, but $R(\alpha)$ big.

i.e. bad for predicting new data.

## Generalization:

Want a decision rule $\hat{\alpha}$ so that $R_{emp}(\hat{\alpha})$ is small, but $R(\hat{\alpha})$ is small.

In practice — cross-validation.

training set $\{(x_i, y_i) : i = 1\ \text{to}\ N\}$
to learn the rule $\bar{\alpha}$

test set $\{(x_j, y_j) : j = 1\ \text{to}\ M\}$
to test the rule $\bar{\alpha}$.

Choose $\hat{\alpha}$ so that $R_{emp}(\hat{\alpha})$ is small on both the training set and test set.
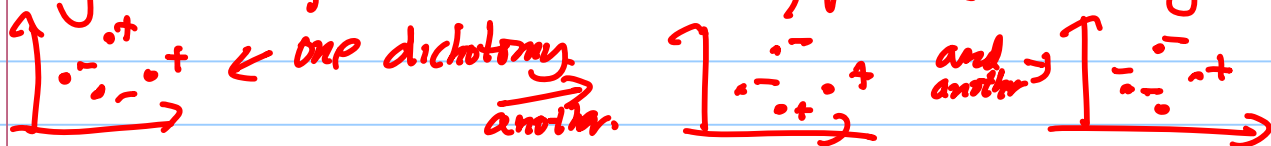How, restrict the possibilities of $\hat{\alpha}$.

What happens if we have an infinite
set of rules? — e.g. the set all seperating planes
$$ax + by + c = 0$$

The Vapnik-Chervonenkis VC dimension gives a
finite measure of the capacity of a hypothesis class $A$.
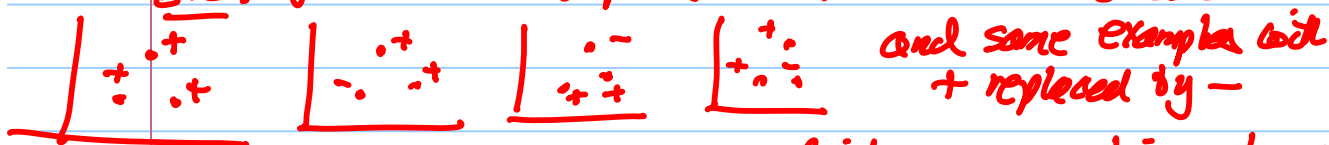
Introduce the concept of shattering.

Suppose we have $n$ data examples (features /attributes) $\{x_i : i = 1...n\}$
in $d$-dim space. With general position assumption (data doesn't
lie on a lower-dimensional subspace).

They are $2^n$ possible dichotomies of the data —
separating the examples into two classes, positive and negative

 ← one dichotomy.  and another → 

A set $A$ of classifiers, shatters $n$ examples in
$d$-dim space if, for all dichotomies of the data,
we can find a classifier in $A$ which classifies the data
correctly.

E.G. If we have 3 datapoints in 2D, there are $2^3 = 8$ dichotomies.

 and some examples with
+ replaced by —

For each dichotomy, we can find a separating plane
which classifies the data perfectly → eg 

Hence, we know that we can classify the data
perfectly before we even look at it.

(16)

The VC-dimension of a hypothesis class $A$ is the maximum number of points that can be shattered. Note: this depends on the dimension of the space.

For seperating hyperplanes, the
VC dimension $= d+1$ . ie. $VC = 3$ for
$\quad\quad\quad\quad\quad\quad\quad$ x dim of space. $\quad\quad\quad$ planes in 2D space.

This concept enables us to proove theorems for hypothesis spaces with finite VC dimension, but infinite number of classifiers (e.g. planes)

For example,
$\quad$ with prob $> 1-\delta$
$$R(\alpha) \leq R_{emp}(\alpha:N) + \sqrt{\frac{h(\log 2N/h) - \log \delta/4}{N}}$$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad$ for all $\alpha \in A$

PAC Theorem $\quad\quad$ where $h$ is the VC dimension of $A$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad$ $N$ is the total amount of data.

Moral: In order to generalize, you have to restrict the complexity (ie. the VC dimension) of the set of classifiers you use by taking into account the amount of data.