# Generative Learning

$$D \sim P(X, Y)$$

r.v.s

generate

$$P(X, Y ; \theta)$$

## MLE (frequentist)

$$\theta = \operatorname{argmax}_{\theta} P(D | \theta)$$

parameter

Which $\theta$ maximizes the
probability of seeing the dataset $D$

## MAP (Bayesian)

$$\theta = \operatorname{arg\,max}_{\theta} P(\theta | D)$$

r.v.

Which $\theta$ is the most likely,
given that we have observed
the data set $D$

## Summary of MAP

- If $n \to +\infty$, $\theta_{MAP} \to \theta_{MLE}$.

- If $n$ is small, the estimate will depends or the prior belief.

- MAP is a great estimator if an accurate prior belief is available.

## True Bayesian Approach

$$P(Y=y | X=x, D) = \int_{\theta} P(Y=y | \theta) P(\theta | D) \, d\theta$$

average out all possible value of $\theta$

Bayes Classifier: return $\underset{\forall y}{\arg\max}\ P(Y=y \mid X=x; \theta)$.

---

## Estimating the conditional probability $P(Y=y \mid X=x)$

- Estimating $P(X=x, Y=y)$ is fine. But why don't we estimate
  $P(Y=y \mid X=x)$ directly.

  Note: If we have enough data, we could estimate $P(X,Y)$
  where we imagine a gigantic die that
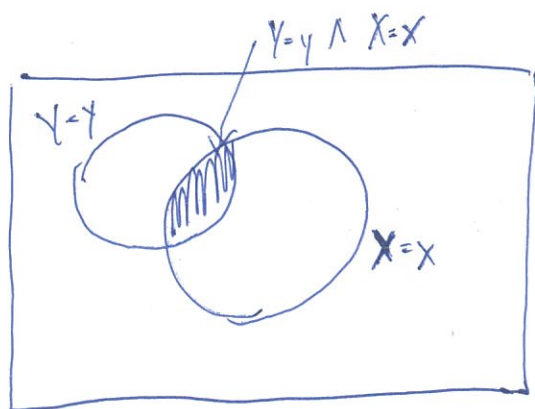  has one side for each possible value of $(X,y)$.



  $P(X=x \wedge Y=y)$ is the probability that one specific
  side coming up.

  EX: By assuming $X, Y$ togethers forms a R.V. that follows
  the binomial distribution, then we have the following by MLE

  $$P(X=x \wedge Y=y) = \sum_{i=1}^{n} \frac{I(X_i=x \wedge y_i=y)}{n}$$

  $$\underset{\text{indicator r.v.}}{I(X_i=x \wedge Y_i=y)} = \begin{cases} 1 & \text{if } x_i=x \text{ and } y_i=y \\ 0, & \text{o.t.w.} \end{cases}$$



  $$P(Y=y \mid X=x) = \frac{P(Y=y \wedge X=x)}{P(X=x)}$$

  $$\Rightarrow P(Y=y \mid X=x) = \frac{\sum_{i=1}^{r} I(X_i=x \wedge y_i=y)}{\cancel{K}}$$

  $$\frac{}{\dfrac{\sum_{i=1}^{n} I(X_i=x)}{\cancel{K}}}$$

  What's the problem?

– In d-dimensional space,

$$P(Y=y \mid \underline{X}=\underline{x}) = P(Y=y \mid [X_1]=x_1, [X_2]=x_2, \ldots, [X_d]=x_d)$$

$$\underline{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{pmatrix} \text{ // feature vector}$$

$$X = \{[X_1], [X_2], \ldots, [X_d]\}, \text{ Each } X_i \text{ is r.v. associated w/ each coordinate}$$

**\*\*\* Problem** (Practical): Estimation by MLE is only good if we have many training vectors w/ same identical features as x. As $d \to +\infty$, this will never happens (e.g. image data)

$\Rightarrow$ If $d \to +\infty$, $P(Y=y, X=x) \to 0$ and $P(X=x) \to 0$

---

## Naive Bayes Classifier:

– By Bayes rule, we have that $P(Y=y \mid X=x) = \dfrac{P(X=x \mid Y=y) \cdot P(Y=y)}{P(X=x)}$

– Again, we are in the world of generative learning

$\Rightarrow$ We have already know how to estimate $P(X=x)$ and $P(Y=y)$

$\Rightarrow$ How about estimating $P(X=x \mid Y=y)$ ??

**Naive Bayes Assumption:** All feature values are independent, given the label

$$P(X=x \mid Y=y) \cdot \prod_{j=1}^{d} P([X_j]=x_j \mid Y=y)$$

- With the assumption, we can derive

$$h(x) = \arg\max_{\forall y} \{ P(Y=y \mid X=x)$$

$$= \arg\max_{\forall y} \frac{P(X=x \mid Y=y)\, P(Y=y)}{P(X=x)}$$

$P(X=x) \leftarrow$ constant (unchange) *(unchange)*

$$= \arg\max_{\forall y} \left( \prod_{j=1}^{d} P([X_j] = x_j \mid Y=y) \right) \cdot P(Y=y)$$

$$= \arg\max_{\forall y} \left[ \sum_{j=1}^{d} \log \left( P([X_j]=x_j \mid Y=y) \right) + \log(P(y)) \right]$$

*easy as one dimension*

---

## Ex: Spam filter by Naive Bayes / text classification

- Each vocabulary is one feature dimension
- We encode each email as a feature vector $x \in \{0, 1\}^{|V|}$
- Each $x_j = 1$ iff the vocabulary $x_j$ appears in the email

$$\boxed{\text{email}} \longrightarrow x = \begin{pmatrix} x_1 \\ \vdots \\ \vdots \\ x_d \end{pmatrix} \begin{matrix} a \\ \vdots \\ \vdots \\ z \end{matrix}$$

- $Y \in \{ \text{SPAM} / \text{NOT-SPAM} \}$

- For a test email $x_t$, we would like to determine

$$P(Y = \text{SPAM} \mid X = x_t) \quad \text{and} \quad P(Y = \text{NOT-SPAM} \mid X = x_t)$$

- $P(Y = SPAM \mid X = x_f) = P(Y = SPAM \mid [X_1] = x_1, [X_2] = x_2, \ldots [X_d] = x_d)$

$$= \left( \prod_{j=1}^{d} \boxed{P([X_j] = x_j \mid Y = SPAM)} \right) \left( P(Y = y) \right)$$

$$\underline{\hspace{5cm}}$$
$$P'(X = x)$$

In next lecture, we explore how to estimate this term

---

Visualization:



$p([x] \mid y = 2)$

$p([x] \mid y = 1)$

$P([X_1] = y = 1)$

$P([X_2] \mid y = 2)$

$\prod_{j=1}^{d} P([X_j] \mid y = 2)$

$\prod_{j=1}^{d} P([X_j] \mid y = 1)$