# Parameter Estimation and Probabilistic Estimators

## Machine Learning

1. Discriminative and generative models.

1. Discriminative and generative models.
2. Parametric distributions.

# Overview

1. Discriminative and generative models.
2. Parametric distributions.
3. Parametric hypothesis classes.

# Overview

1. Discriminative and generative models.
2. Parametric distributions.
3. Parametric hypothesis classes.
4. Bayes rule.

# Overview

1. Discriminative and generative models.
2. Parametric distributions.
3. Parametric hypothesis classes.
4. Bayes rule.
5. Bayesian parameter estimation.

# Overview

1. Discriminative and generative models.
2. Parametric distributions.
3. Parametric hypothesis classes.
4. Bayes rule.
5. Bayesian parameter estimation.
6. Naive Bayes.

# Discriminative and Generative Models

# Discriminative and Generative Models

Two main types of models in supervised learning:

# Discriminative and Generative Models

Two main types of models in supervised learning:

### Definition (Discriminative Model)

A **discriminative model** is one that learns and models $p_{Y|X}(y \mid x)$. In other words, it learns the *conditional* distribution and is useful for classification/regression.

# Discriminative and Generative Models

Two main types of models in supervised learning:

## Definition (Discriminative Model)

A **discriminative model** is one that learns and models $p_{Y|X}(y \mid x)$. In other words, it learns the *conditional* distribution and is useful for classification/regression.

## Definition (Generative Model)

A **generative model** is one that learns and models $p_{X,Y}(x, y)$. In other words, it learns the *joint* distribution and, on top of classification/regression, can be used for other tasks such as generating plausible data.

# Discriminative and Generative Models

Two main types of models in supervised learning:

### Definition (Discriminative Model)

A **discriminative model** is one that learns and models $p_{Y|X}(y \mid x)$. In other words, it learns the *conditional* distribution and is useful for classification/regression.

### Definition (Generative Model)

A **generative model** is one that learns and models $p_{X,Y}(x, y)$. In other words, it learns the *joint* distribution and, on top of classification/regression, can be used for other tasks such as generating plausible data.

Generative models are more versatile; discriminative models are more effective (Ng, Jordan).

## Parametric Distributions

Sometimes we want to model our data as from a distribution $\mu_\lambda$ in a family of distributions

## Parametric Distributions

Sometimes we want to model our data as from a distribution $\mu_\lambda$ in a family of distributions

$$\mathcal{M} = \{\mu_\lambda : \lambda \in \Lambda\}.$$

## Parametric Distributions

Sometimes we want to model our data as from a distribution $\mu_\lambda$ in a family of distributions

$$\mathcal{M} = \{\mu_\lambda : \lambda \in \Lambda\}.$$

We could know this by prior knowledge,

## Parametric Distributions

Sometimes we want to model our data as from a distribution $\mu_\lambda$ in a family of distributions

$$\mathcal{M} = \{\mu_\lambda : \lambda \in \Lambda\}.$$

We could know this by prior knowledge, *or* we could just be guessing and hoping the true distribution is close enough to one of these distributions.

## Parametric Distributions

Sometimes we want to model our data as from a distribution $\mu_\lambda$ in a family of distributions

$$\mathcal{M} = \{\mu_\lambda : \lambda \in \Lambda\}.$$

We could know this by prior knowledge, *or* we could just be guessing and hoping the true distribution is close enough to one of these distributions. Here $\lambda$ is the **parameter**.

## Parametric Distributions

Sometimes we want to model our data as from a distribution $\mu_\lambda$ in a family of distributions

$$\mathcal{M} = \{\mu_\lambda : \lambda \in \Lambda\}.$$

We could know this by prior knowledge, *or* we could just be guessing and hoping the true distribution is close enough to one of these distributions. Here $\lambda$ is the **parameter**. Correspondingly, for unsupervised problems we can consider $\mu_{\lambda;X}$ coming from the family

## Parametric Distributions

Sometimes we want to model our data as from a distribution $\mu_\lambda$ in a family of distributions

$$\mathcal{M} = \{\mu_\lambda : \lambda \in \Lambda\}.$$

We could know this by prior knowledge, *or* we could just be guessing and hoping the true distribution is close enough to one of these distributions. Here $\lambda$ is the **parameter**. Correspondingly, for unsupervised problems we can consider $\mu_{\lambda;X}$ coming from the family

$$\mathcal{M}_X = \{\mu_{\lambda;X} : \lambda \in \Lambda\}.$$

# Parametric Distributions

Sometimes we want to model our data as from a distribution $\mu_\lambda$ in a family of distributions

$$\mathcal{M} = \{\mu_\lambda : \lambda \in \Lambda\}.$$

We could know this by prior knowledge, *or* we could just be guessing and hoping the true distribution is close enough to one of these distributions. Here $\lambda$ is the **parameter**. Correspondingly, for unsupervised problems we can consider $\mu_{\lambda;X}$ coming from the family

$$\mathcal{M}_X = \{\mu_{\lambda;X} : \lambda \in \Lambda\}.$$

## Example (Normal Family)

If we suppose our data is normally distributed, we can let $\lambda = (\theta, \Sigma)$, so

$$\mathcal{M}_X = \left\{ E \mapsto \mathsf{P}[\mathcal{N}(\theta, \Sigma) \in E] : \theta \in \mathbb{R}^d, \Sigma \in \mathbb{R}^{d \times d} \right\}.$$

Sometimes our algorithm gives a hypothesis class of the form

Sometimes our algorithm gives a hypothesis class of the form

$$\mathcal{H} = \{h_\theta : \theta \in \Theta\} \, .$$

Sometimes our algorithm gives a hypothesis class of the form

$$\mathcal{H} = \{h_\theta : \theta \in \Theta\}.$$

This is entirely dependent on the algorithm we choose.

## Parametric Hypotheses

Sometimes our algorithm gives a hypothesis class of the form

$$\mathcal{H} = \{h_\theta \colon \theta \in \Theta\}.$$

This is entirely dependent on the algorithm we choose.

### Example (Support Vector Classification)

The popular support vector machine algorithm has hypotheses of the form

$$\mathcal{H} = \left\{x \mapsto \text{sign}(\langle \theta, x \rangle) \colon \theta \in \mathbb{R}^d\right\}.$$

# Bayes' Rule

## Theorem (Bayes' Rule)

# Bayes' Rule

## Theorem (Bayes' Rule)

*If $A$ and $B$ are events, then*

$$\mathsf{P}[B \mid A] = \frac{\mathsf{P}[A \mid B]\,\mathsf{P}[B]}{\mathsf{P}[A]}.$$

# Bayes' Rule

## Theorem (Bayes' Rule)

*If $A$ and $B$ are events, then*

$$P[B \mid A] = \frac{P[A \mid B] P[B]}{P[A]}.$$

*If $x$ and $y$ are random variables with density $p_X$ and $p_Y$, then*

$$p_{X|Y}(u \mid v) = \frac{p_{Y|X}(v \mid u) p_X(u)}{p_Y(v)}.$$

# Bayes' Rule

## Theorem (Bayes' Rule)

*If $A$ and $B$ are events, then*

$$P[B \mid A] = \frac{P[A \mid B] \, P[B]}{P[A]}.$$

*If $x$ and $y$ are random variables with density $p_X$ and $p_Y$, then*

$$p_{X|Y}(u \mid v) = \frac{p_{Y|X}(v \mid u) \, p_X(u)}{p_Y(v)}.$$

Proof is not very complicated:

# Bayes' Rule

## Theorem (Bayes' Rule)

*If $A$ and $B$ are events, then*

$$P[B \mid A] = \frac{P[A \mid B] \, P[B]}{P[A]}.$$

*If $x$ and $y$ are random variables with density $p_X$ and $p_Y$, then*

$$p_{X|Y}(u \mid v) = \frac{p_{Y|X}(v \mid u) \, p_X(u)}{p_Y(v)}.$$

Proof is not very complicated:

$$P[B \mid A] \, P[A] = P[A \cap B] = P[A \mid B] \, P[B].$$

Let's say we have a distribution class

Let's say we have a distribution class

$$\mathcal{M} = \{\mu_\lambda : \lambda \in \Lambda\}.$$

Let's say we have a distribution class

$$\mathcal{M} = \{\mu_\lambda : \lambda \in \Lambda\}.$$

Data distributed according to $\mu_{\lambda_\star}$

Let's say we have a distribution class

$$\mathcal{M} = \{\mu_\lambda : \lambda \in \Lambda\}.$$

Data distributed according to $\mu_{\lambda_\star}$; want to estimate $\lambda_\star$ using the **estimator** $\widehat{\lambda} = \widehat{\lambda}(s)$.

Simplest, most intuitive parameter estimation algorithm.

# MLE

Simplest, most intuitive parameter estimation algorithm.

## MLE Algorithm

- **Input:** a distribution class $\mathcal{M}$, a sample $s$.
- **Output:** an estimate $\widehat{\lambda}_{\mathrm{MLE}}(s)$ for $\lambda_\star$.
- **Algorithm:** *Take the $\lambda$ making your data most likely.*

$$\widehat{\lambda}_{\mathrm{MLE}}(s) \in \underset{\lambda \in \Lambda}{\operatorname{argmax}}\; p_S(s \mid \lambda).$$

Logarithm is monotonic, so we have some equivalent forms:

## MLE

Logarithm is monotonic, so we have some equivalent forms:

$$\widehat{\lambda}_{\mathrm{MLE}}(s)$$

# MLE

Logarithm is monotonic, so we have some equivalent forms:

$$\widehat{\lambda}_{\mathrm{MLE}}(s) \in \operatorname*{argmax}_{\lambda \in \Lambda} p_S(s \mid \lambda)$$

# MLE

Logarithm is monotonic, so we have some equivalent forms:

$$\widehat{\lambda}_{\mathrm{MLE}}(s) \in \underset{\lambda \in \Lambda}{\operatorname{argmax}} \, p_S(s \mid \lambda)$$

$$= \underset{\lambda \in \Lambda}{\operatorname{argmax}} \prod_{i=1}^{n} p_Z(z_i \mid \lambda)$$

# MLE

Logarithm is monotonic, so we have some equivalent forms:

$$\widehat{\lambda}_{\mathrm{MLE}}(s) \in \underset{\lambda \in \Lambda}{\operatorname{argmax}}\ p_S(s \mid \lambda)$$

$$= \underset{\lambda \in \Lambda}{\operatorname{argmax}} \prod_{i=1}^{n} p_Z(z_i \mid \lambda)$$

$$= \underset{\lambda \in \Lambda}{\operatorname{argmax}} \sum_{i=1}^{n} \log(p_Z(z_i \mid \lambda)).$$

# MLE

Logarithm is monotonic, so we have some equivalent forms:

$$\widehat{\lambda}_{\text{MLE}}(s) \in \underset{\lambda \in \Lambda}{\arg\max}\ p_S(s \mid \lambda)$$

$$= \underset{\lambda \in \Lambda}{\arg\max}\ \prod_{i=1}^{n} p_Z(z_i \mid \lambda)$$

$$= \underset{\lambda \in \Lambda}{\arg\max}\ \sum_{i=1}^{n} \log(p_Z(z_i \mid \lambda)).$$

*In supervised learning, we should factor $p_Z = p_{X,Y}$ into $p_{X|Y}\, p_Y$ or $p_{Y|X}\, p_X$, whichever is easier or more efficient to compute.*

## MLE

Logarithm is monotonic, so we have some equivalent forms:

$$\widehat{\lambda}_{\mathrm{MLE}}(s) \in \operatorname*{argmax}_{\lambda \in \Lambda} p_S(s \mid \lambda)$$

$$= \operatorname*{argmax}_{\lambda \in \Lambda} \prod_{i=1}^{n} p_Z(z_i \mid \lambda)$$

$$= \operatorname*{argmax}_{\lambda \in \Lambda} \sum_{i=1}^{n} \log(p_Z(z_i \mid \lambda)).$$

*In supervised learning, we should factor $p_Z = p_{X,Y}$ into $p_{X|Y} \, p_Y$ or $p_{Y|X} \, p_X$, whichever is easier or more efficient to compute.*

*In unsupervised learning, we can just treat $p_Z = p_X$ and compute the density directly.*

Above solution for MLE was straightforward and simple, yet it works pretty well!

Above solution for MLE was straightforward and simple, yet it works pretty well!
One thing it's missing is that it doesn't let us add prior information about what we *think* $\lambda_\star$ is.

Above solution for MLE was straightforward and simple, yet it works pretty well!
One thing it's missing is that it doesn't let us add prior information about what we *think* $\lambda_\star$ is.
**Solution:** treat $\lambda$ like a *random variable*, and give it a distribution $\mu_\Lambda$ and density $p_\Lambda$ which reflects prior information.

More advanced parameter estimation algorithm.

# MAP

More advanced parameter estimation algorithm.

## MAP Algorithm

- **Input:** a distribution class $\mathcal{M}$, a sample $s$, a prior $p_\Lambda$ on the value of $\lambda_\star$.
- **Output:** an estimate $\widehat{\lambda}_{\mathrm{MAP}}(s)$ for $\lambda_\star$.
- **Algorithm:** *Take the $\lambda$ which is most likely, given the data and your prior.*

$$\widehat{\lambda}_{\mathrm{MAP}}(s) \in \underset{\lambda \in \Lambda}{\mathrm{argmax}}\, p_{\Lambda|S}(\lambda \mid s).$$

# MAP

We can use monotonicity of the logarithm again:

# MAP

We can use monotonicity of the logarithm again:

$$\widehat{\lambda}_{\text{MAP}}(s)$$

# MAP

We can use monotonicity of the logarithm again:

$$\widehat{\lambda}_{\mathrm{MAP}}(s) = \operatorname*{argmax}_{\lambda \in \Lambda} p_{\Lambda \mid S}(\lambda \mid s)$$

## MAP

We can use monotonicity of the logarithm again:

$$\widehat{\lambda}_{\mathrm{MAP}}(s) = \underset{\lambda \in \Lambda}{\mathrm{argmax}}\, p_{\Lambda|S}(\lambda \mid s)$$
$$= \underset{\lambda \in \Lambda}{\mathrm{argmax}}\, \frac{p_{S|\Lambda}(s \mid \lambda)\, p_{\Lambda}(\lambda)}{p_S(s)}$$

# MAP

We can use monotonicity of the logarithm again:

$$
\begin{aligned}
\widehat{\lambda}_{\mathrm{MAP}}(s) &= \underset{\lambda \in \Lambda}{\operatorname{argmax}}\, p_{\Lambda|S}(\lambda \mid s) \\
&= \underset{\lambda \in \Lambda}{\operatorname{argmax}}\, \frac{p_{S|\Lambda}(s \mid \lambda)\, p_{\Lambda}(\lambda)}{p_S(s)} \\
&= \underset{\lambda \in \Lambda}{\operatorname{argmax}}\, p_{S|\Lambda}(s \mid \lambda)\, p_{\Lambda}(\lambda)
\end{aligned}
$$

## MAP

We can use monotonicity of the logarithm again:

$$
\begin{aligned}
\widehat{\lambda}_{\mathrm{MAP}}(s) &= \underset{\lambda \in \Lambda}{\operatorname{argmax}} \; p_{\Lambda|S}(\lambda \mid s) \\
&= \underset{\lambda \in \Lambda}{\operatorname{argmax}} \; \frac{p_{S|\Lambda}(s \mid \lambda) p_{\Lambda}(\lambda)}{p_S(s)} \\
&= \underset{\lambda \in \Lambda}{\operatorname{argmax}} \; p_{S|\Lambda}(s \mid \lambda) p_{\Lambda}(\lambda) \\
&= \underset{\lambda \in \Lambda}{\operatorname{argmax}} \; p_{\Lambda}(\lambda) \prod_{i=1}^{n} p_{Z|\Lambda}(z_i \mid \lambda)
\end{aligned}
$$

# MAP

We can use monotonicity of the logarithm again:

$$
\begin{aligned}
\widehat{\lambda}_{\mathrm{MAP}}(s) &= \operatorname*{argmax}_{\lambda \in \Lambda} p_{\Lambda|S}(\lambda \mid s) \\
&= \operatorname*{argmax}_{\lambda \in \Lambda} \frac{p_{S|\Lambda}(s \mid \lambda) p_{\Lambda}(\lambda)}{p_S(s)} \\
&= \operatorname*{argmax}_{\lambda \in \Lambda} p_{S|\Lambda}(s \mid \lambda) p_{\Lambda}(\lambda) \\
&= \operatorname*{argmax}_{\lambda \in \Lambda} p_{\Lambda}(\lambda) \prod_{i=1}^{n} p_{Z|\Lambda}(z_i \mid \lambda) \\
&= \operatorname*{argmax}_{\lambda \in \Lambda} \left( \log(p_{\Lambda}(\lambda)) + \sum_{i=1}^{n} \log\big(p_{Z|\Lambda}(z_i \mid \lambda)\big) \right).
\end{aligned}
$$

## MAP

We can use monotonicity of the logarithm again:

$$
\begin{aligned}
\widehat{\lambda}_{\mathrm{MAP}}(s) &= \operatorname*{argmax}_{\lambda \in \Lambda} p_{\Lambda|S}(\lambda \mid s) \\
&= \operatorname*{argmax}_{\lambda \in \Lambda} \frac{p_{S|\Lambda}(s \mid \lambda) p_{\Lambda}(\lambda)}{p_S(s)} \\
&= \operatorname*{argmax}_{\lambda \in \Lambda} p_{S|\Lambda}(s \mid \lambda) p_{\Lambda}(\lambda) \\
&= \operatorname*{argmax}_{\lambda \in \Lambda} p_{\Lambda}(\lambda) \prod_{i=1}^{n} p_{Z|\Lambda}(z_i \mid \lambda) \\
&= \operatorname*{argmax}_{\lambda \in \Lambda} \left( \log(p_{\Lambda}(\lambda)) + \sum_{i=1}^{n} \log\big(p_{Z|\Lambda}(z_i \mid \lambda)\big) \right).
\end{aligned}
$$

*We can do the same factorization of $p_{Z|\Lambda}$ as MLE.*

# Naive Bayes

Generative supervised learning algorithm for classification using parameter estimation.

# Naive Bayes

Generative supervised learning algorithm for classification using parameter estimation.

## Naive Bayes Algorithm

- **Input:** A known or learned parameter $\lambda$, a new data point $x \in \mathcal{X}$.
- **Output:** A predicted output $\widehat{y}(x)$ corresponding to $x$.
- **Algorithm:** *Maximize the posterior probability:*

$$\widehat{y}(x) \in \underset{y \in \mathcal{Y}}{\operatorname{argmax}} \, p_{Y|X}(y \mid x, \lambda).$$

# Naive Bayes

Use the same tricks as MAP to compute.

# Naive Bayes

Use the same tricks as MAP to compute.

$$\widehat{y}(x)$$

# Naive Bayes

Use the same tricks as MAP to compute.

$$\widehat{y}(x) \in \operatorname*{argmax}_{y \in \mathcal{Y}} p_{Y|X}(y \mid x, \lambda)$$

# Naive Bayes

Use the same tricks as MAP to compute.

$$\widehat{y}(x) \in \underset{y \in \mathcal{Y}}{\operatorname{argmax}} \, p_{Y|X}(y \mid x, \lambda)$$

$$= \underset{y \in \mathcal{Y}}{\operatorname{argmax}} \, \frac{p_{X,Y}(x, y \mid \lambda)}{p_X(x)}$$

# Naive Bayes

Use the same tricks as MAP to compute.

$$\widehat{y}(x) \in \operatorname*{argmax}_{y \in \mathcal{Y}} p_{Y|X}(y \mid x, \lambda)$$

$$= \operatorname*{argmax}_{y \in \mathcal{Y}} \frac{p_{X,Y}(x, y \mid \lambda)}{p_X(x)}$$

$$= \operatorname*{argmax}_{y \in \mathcal{Y}} p_{X,Y}(x, y \mid \lambda)$$

# Naive Bayes

Use the same tricks as MAP to compute.

$$\widehat{y}(x) \in \operatorname*{argmax}_{y \in \mathcal{Y}} p_{Y|X}(y \mid x, \lambda)$$

$$= \operatorname*{argmax}_{y \in \mathcal{Y}} \frac{p_{X,Y}(x, y \mid \lambda)}{p_X(x)}$$

$$= \operatorname*{argmax}_{y \in \mathcal{Y}} p_{X,Y}(x, y \mid \lambda)$$

$$= \operatorname*{argmax}_{y \in \mathcal{Y}} p_{X|Y}(x \mid y, \lambda) p_Y(y \mid \lambda)$$

# Naive Bayes

Use the same tricks as MAP to compute.

$$
\begin{aligned}
\widehat{y}(x) &\in \underset{y \in \mathcal{Y}}{\operatorname{argmax}}\, p_{Y|X}(y \mid x, \lambda) \\
&= \underset{y \in \mathcal{Y}}{\operatorname{argmax}}\, \frac{p_{X,Y}(x, y \mid \lambda)}{p_X(x)} \\
&= \underset{y \in \mathcal{Y}}{\operatorname{argmax}}\, p_{X,Y}(x, y \mid \lambda) \\
&= \underset{y \in \mathcal{Y}}{\operatorname{argmax}}\, p_{X|Y}(x \mid y, \lambda)\, p_Y(y \mid \lambda) \\
&= \underset{y \in \mathcal{Y}}{\operatorname{argmax}}\, \left( \log\big(p_{X|Y}(x \mid y, \lambda)\big) + \log(p_Y(y \mid \lambda)) \right).
\end{aligned}
$$

## Naive Bayes

Use the same tricks as MAP to compute.

$$
\begin{aligned}
\widehat{y}(x) &\in \operatorname*{argmax}_{y \in \mathcal{Y}} p_{Y|X}(y \mid x, \lambda) \\
&= \operatorname*{argmax}_{y \in \mathcal{Y}} \frac{p_{X,Y}(x, y \mid \lambda)}{p_X(x)} \\
&= \operatorname*{argmax}_{y \in \mathcal{Y}} p_{X,Y}(x, y \mid \lambda) \\
&= \operatorname*{argmax}_{y \in \mathcal{Y}} p_{X|Y}(x \mid y, \lambda) p_Y(y \mid \lambda) \\
&= \operatorname*{argmax}_{y \in \mathcal{Y}} \left( \log\left(p_{X|Y}(x \mid y, \lambda)\right) + \log(p_Y(y \mid \lambda)) \right).
\end{aligned}
$$

If each class's points are drawn from a separate distribution, then $\lambda$ contains the parameters for each distribution, and $Y$ gives the distribution to evaluate.

## Naive Bayes

Use the same tricks as MAP to compute.

$$
\begin{aligned}
\widehat{y}(x) &\in \operatorname*{argmax}_{y \in \mathcal{Y}} p_{Y|X}(y \mid x, \lambda) \\
&= \operatorname*{argmax}_{y \in \mathcal{Y}} \frac{p_{X,Y}(x, y \mid \lambda)}{p_X(x)} \\
&= \operatorname*{argmax}_{y \in \mathcal{Y}} p_{X,Y}(x, y \mid \lambda) \\
&= \operatorname*{argmax}_{y \in \mathcal{Y}} p_{X|Y}(x \mid y, \lambda) p_Y(y \mid \lambda) \\
&= \operatorname*{argmax}_{y \in \mathcal{Y}} \left( \log\!\left( p_{X|Y}(x \mid y, \lambda) \right) + \log( p_Y(y \mid \lambda)) \right).
\end{aligned}
$$

If each class's points are drawn from a separate distribution, then $\lambda$ contains the parameters for each distribution, and $Y$ gives the distribution to evaluate. Thus, $\lambda$ lets us evaluate $p_{X|Y}$.

## Naive Bayes

Use the same tricks as MAP to compute.

$$
\begin{aligned}
\widehat{y}(x) &\in \underset{y \in \mathcal{Y}}{\operatorname{argmax}}\, p_{Y|X}(y \mid x, \lambda) \\
&= \underset{y \in \mathcal{Y}}{\operatorname{argmax}}\, \frac{p_{X,Y}(x, y \mid \lambda)}{p_X(x)} \\
&= \underset{y \in \mathcal{Y}}{\operatorname{argmax}}\, p_{X,Y}(x, y \mid \lambda) \\
&= \underset{y \in \mathcal{Y}}{\operatorname{argmax}}\, p_{X|Y}(x \mid y, \lambda)\, p_Y(y \mid \lambda) \\
&= \underset{y \in \mathcal{Y}}{\operatorname{argmax}} \left( \log\left(p_{X|Y}(x \mid y, \lambda)\right) + \log\left(p_Y(y \mid \lambda)\right) \right).
\end{aligned}
$$

If each class's points are drawn from a separate distribution, then $\lambda$ contains the parameters for each distribution, and $Y$ gives the distribution to evaluate. Thus, $\lambda$ lets us evaluate $p_{X|Y}$. We estimate $p_Y$ using data by just counting proportion in each class.

## Naive Bayes

Use the same tricks as MAP to compute.

$$\begin{aligned}
\widehat{y}(x) &\in \operatorname*{argmax}_{y \in \mathcal{Y}} p_{Y|X}(y \mid x, \lambda) \\
&= \operatorname*{argmax}_{y \in \mathcal{Y}} \frac{p_{X,Y}(x, y \mid \lambda)}{p_X(x)} \\
&= \operatorname*{argmax}_{y \in \mathcal{Y}} p_{X,Y}(x, y \mid \lambda) \\
&= \operatorname*{argmax}_{y \in \mathcal{Y}} p_{X|Y}(x \mid y, \lambda) p_Y(y \mid \lambda) \\
&= \operatorname*{argmax}_{y \in \mathcal{Y}} \left( \log\left( p_{X|Y}(x \mid y, \lambda) \right) + \log( p_Y(y \mid \lambda)) \right).
\end{aligned}$$

If each class's points are drawn from a separate distribution, then $\lambda$ contains the parameters for each distribution, and $Y$ gives the distribution to evaluate. Thus, $\lambda$ lets us evaluate $p_{X|Y}$. We estimate $p_Y$ using data by just counting proportion in each class.

Quick and easy classifier!

Machine learning algorithms usually have a **training phase** and **inference phase**.

Machine learning algorithms usually have a **training phase** and **inference phase**.

- The training phase is where the model is constructed and updated.

Machine learning algorithms usually have a **training phase** and **inference phase**.

- The training phase is where the model is constructed and updated.
- The inference phase is where the model is applied to data it hasn't seen before.

Machine learning algorithms usually have a **training phase** and **inference phase**.

- The training phase is where the model is constructed and updated.
- The inference phase is where the model is applied to data it hasn't seen before.

In SKLearn these are encapsulated by the APIs model.fit and model.predict.

## Naive Bayes Algorithm

**procedure** TRAIN($s_{\text{train}}$)
    $\widehat{\lambda} \leftarrow \text{Estimate}(s_{\text{train}})$          ▷ Parameter estimation via MLE.
    **for** $(x_i, y_i) = z_i \in s_{\text{train}}$ **do**
        $\widehat{p}_Y(y_i \mid \lambda) \leftarrow \widehat{p}_Y(y_i \mid \lambda) + \frac{1}{|s_{\text{train}}|}$          ▷ Estimates $p_Y$ by $\widehat{p}_Y$.

## Naive Bayes Algorithm

**procedure** TRAIN($s_{\text{train}}$)
    $\widehat{\lambda} \leftarrow \texttt{Estimate}(s_{\text{train}})$                                                           $\triangleright$ Parameter estimation via MLE.
    **for** $(x_i, y_i) = z_i \in s_{\text{train}}$ **do**
        $\widehat{p}_Y(y_i \mid \lambda) \leftarrow \widehat{p}_Y(y_i \mid \lambda) + \frac{1}{|s_{\text{train}}|}$                                       $\triangleright$ Estimates $p_Y$ by $\widehat{p}_Y$.

**procedure** INFERENCE($x$)
    **return** $\operatorname{argmax}_{y \in \mathcal{Y}} \left( \underbrace{\log\left( p_{X|Y}(x \mid y, \widehat{\lambda}) \right)}_{\text{closed form}} + \underbrace{\log\left( \widehat{p}_Y(y \mid \widehat{\lambda}) \right)}_{\text{estimated in training}} \right)$      $\triangleright$ Iterates through all
$y \in \mathcal{Y}$.

# Recap

Main concepts:

- **Generative models**: learn $p_{X,Y}(x, y)$
- **Discriminative models**: learn $p_{Y|X}(y \mid x)$
- **MLE**:

$$\widehat{\lambda}_{\text{MLE}}(s) \in \underset{\lambda \in \Lambda}{\operatorname{argmax}} \; p_S(s \mid \lambda)$$

- **MAP**:

$$\widehat{\lambda}_{\text{MAP}}(s) \in \underset{\lambda \in \Lambda}{\operatorname{argmax}} \; p_{\Lambda|S}(\lambda \mid s)$$

- **Naive Bayes**:

$$\widehat{y}(x) \in \underset{y \in \mathcal{Y}}{\operatorname{argmax}} \; p_{Y|X}(y \mid x, \lambda)$$