

Probabilistic Models for Classification: Discriminative Classification

Piyush Rai

Probabilistic Machine Learning (CS772A)

Aug 19, 2017

Recap

Generative Classification

- A two-step procedure to learn a probabilistic classification model $p(y|\mathbf{x})$

Generative Classification

- A two-step procedure to learn a probabilistic classification model $p(y|\mathbf{x})$
 - ① Learn **class-conditional distribution** $p(\mathbf{x}|y)$ and **class-prior distribution** $p(y)$ using training data

Generative Classification

- A two-step procedure to learn a probabilistic classification model $p(y|\mathbf{x})$
 - ① Learn **class-conditional distribution** $p(\mathbf{x}|y)$ and **class-prior distribution** $p(y)$ using training data
 - We usually **assume the form** of $p(\mathbf{x}|y)$ and $p(y)$, but the parameters are unknowns

Generative Classification

- A two-step procedure to learn a probabilistic classification model $p(y|\mathbf{x})$
 - ① Learn **class-conditional distribution** $p(\mathbf{x}|y)$ and **class-prior distribution** $p(y)$ using training data
 - We usually **assume the form** of $p(\mathbf{x}|y)$ and $p(y)$, but the parameters are unknowns
 - Parameter estimation can be done via MLE or MAP or fully Bayesian inference

Generative Classification

- A two-step procedure to learn a probabilistic classification model $p(y|\mathbf{x})$
 - ① Learn **class-conditional distribution** $p(\mathbf{x}|y)$ and **class-prior distribution** $p(y)$ using training data
 - We usually **assume the form** of $p(\mathbf{x}|y)$ and $p(y)$, but the parameters are unknowns
 - Parameter estimation can be done via MLE or MAP or fully Bayesian inference
 - ② Apply Bayes rule to predict the posterior **probability** of each class for any (test) input \mathbf{x}_*

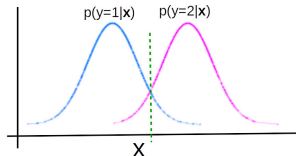
$$p(y_* = k|\mathbf{x}_*) = \frac{p(y_* = k)p(\mathbf{x}_*|y_* = k)}{p(\mathbf{x}_*)}$$

Generative Classification

- A two-step procedure to learn a probabilistic classification model $p(y|\mathbf{x})$
 - 1 Learn **class-conditional distribution** $p(\mathbf{x}|y)$ and **class-prior distribution** $p(y)$ using training data
 - We usually **assume the form** of $p(\mathbf{x}|y)$ and $p(y)$, but the parameters are unknowns
 - Parameter estimation can be done via MLE or MAP or fully Bayesian inference
 - 2 Apply Bayes rule to predict the posterior **probability** of each class for any (test) input \mathbf{x}_*

$$p(y_* = k|\mathbf{x}_*) = \frac{p(y_* = k)p(\mathbf{x}_*|y_* = k)}{p(\mathbf{x}_*)}$$

- The predicted class for the input \mathbf{x} will be the one that has the largest posterior probability



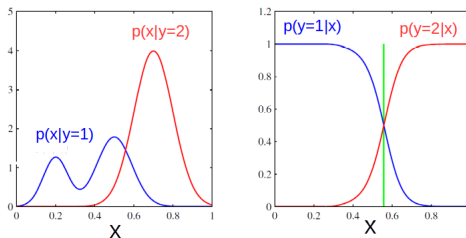
Generative Classification

- Generative classification requires modeling the input distribution (for each class)

Generative Classification

- Generative classification requires modeling the input distribution (for each class)
- In some cases, this extensive modeling effort may not be worth it (e.g., the case shown below)

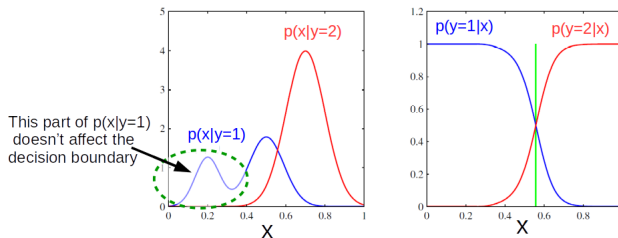
[Assuming $p(y=1) = p(y=2)$]



Generative Classification

- Generative classification requires modeling the input distribution (for each class)
- In some cases, this extensive modeling effort may not be worth it (e.g., the case shown below)

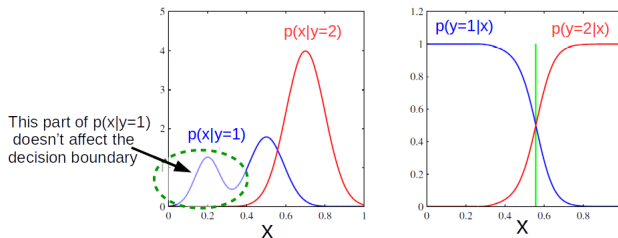
[Assuming $p(y=1) = p(y=2)$]



Generative Classification

- Generative classification requires modeling the input distribution (for each class)
- In some cases, this extensive modeling effort may not be worth it (e.g., the case shown below)

[Assuming $p(y=1) = p(y=2)$]

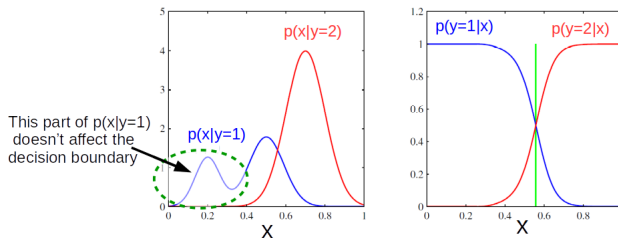


- A better approach would be to directly model $p(y|x)$, i.e., the distribution that makes the decision

Generative Classification

- Generative classification requires modeling the input distribution (for each class)
- In some cases, this extensive modeling effort may not be worth it (e.g., the case shown below)

[Assuming $p(y=1) = p(y=2)$]

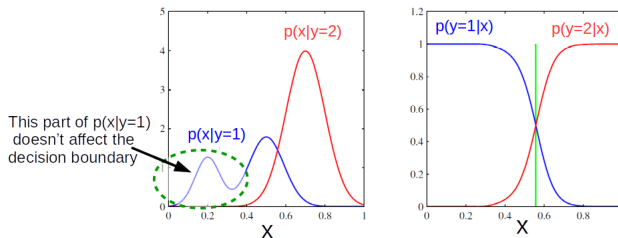


- A better approach would be to directly model $p(y|x)$, i.e., the distribution that makes the decision
 - Such an approach is called **Discriminative Classification** (today's topic)

Generative Classification

- Generative classification requires modeling the input distribution (for each class)
- In some cases, this extensive modeling effort may not be worth it (e.g., the case shown below)

[Assuming $p(y=1) = p(y=2)$]



- A better approach would be to directly model $p(y|x)$, i.e., the distribution that makes the decision
 - Such an approach is called **Discriminative Classification** (today's topic)
- PS: The generative approach is still appealing for many other reasons (as discussed in last lecture)

Probabilistic Models for Discriminative Classification

Probabilistic Models for Discriminative Classification

Note: Many **non-probabilistic** classification models can be termed as discriminative classifiers (e.g., **Support Vector Machine** or SVM, which directly learns a separator between classes)

(Probabilistic) Discriminative Classification

- Models $p(y|\mathbf{x})$ directly using an appropriate discrete distribution

(Probabilistic) Discriminative Classification

- Models $p(y|\mathbf{x})$ directly using an appropriate discrete distribution
- Parameters of the distribution $p(y|\mathbf{x})$ defined by a linear/nonlinear function of \mathbf{x}

(Probabilistic) Discriminative Classification

- Models $p(y|\mathbf{x})$ directly using an appropriate **discrete** distribution
- Parameters of the distribution $p(y|\mathbf{x})$ defined by a **linear/nonlinear function** of \mathbf{x}
 - A discriminative model only needs to learn these parameters!

(Probabilistic) Discriminative Classification

- Models $p(y|\mathbf{x})$ directly using an appropriate **discrete** distribution
- Parameters of the distribution $p(y|\mathbf{x})$ defined by a **linear/nonlinear function** of \mathbf{x}
 - A discriminative model only needs to learn these parameters!
- For binary classification, $p(y|\mathbf{x})$ would be a **Bernoulli**

$$p(y|\mathbf{x}) = p(y|\mathbf{x}, f) = \text{Bernoulli}[f(\mathbf{x})]$$

(Probabilistic) Discriminative Classification

- Models $p(y|\mathbf{x})$ directly using an appropriate **discrete** distribution
- Parameters of the distribution $p(y|\mathbf{x})$ defined by a **linear/nonlinear function** of \mathbf{x}
 - A discriminative model only needs to learn these parameters!
- For binary classification, $p(y|\mathbf{x})$ would be a **Bernoulli**

$$p(y|\mathbf{x}) = p(y|\mathbf{x}, f) = \text{Bernoulli}[f(\mathbf{x})]$$

where $f()$ is a function (**defined by some params**) that maps \mathbf{x} to a probability $\mu = p(y = 1|\mathbf{x})$

(Probabilistic) Discriminative Classification

- Models $p(y|\mathbf{x})$ directly using an appropriate **discrete** distribution
- Parameters of the distribution $p(y|\mathbf{x})$ defined by a **linear/nonlinear function** of \mathbf{x}
 - A discriminative model only needs to learn these parameters!
- For binary classification, $p(y|\mathbf{x})$ would be a **Bernoulli**

$$p(y|\mathbf{x}) = p(y|\mathbf{x}, f) = \text{Bernoulli}[f(\mathbf{x})]$$

where $f()$ is a function (**defined by some params**) that maps \mathbf{x} to a probability $\mu = p(y = 1|\mathbf{x})$

- For $K > 2$ class multiclass classification, $p(y|\mathbf{x})$ would be a **multinoulli**

(Probabilistic) Discriminative Classification

- Models $p(y|\mathbf{x})$ directly using an appropriate **discrete** distribution
- Parameters of the distribution $p(y|\mathbf{x})$ defined by a **linear/nonlinear function** of \mathbf{x}
 - A discriminative model only needs to learn these parameters!
- For binary classification, $p(y|\mathbf{x})$ would be a **Bernoulli**

$$p(y|\mathbf{x}) = p(y|\mathbf{x}, f) = \text{Bernoulli}[f(\mathbf{x})]$$

where $f()$ is a function (**defined by some params**) that maps \mathbf{x} to a probability $\mu = p(y = 1|\mathbf{x})$

- For $K > 2$ class multiclass classification, $p(y|\mathbf{x})$ would be a **multinoulli**

$$p(y|\mathbf{x}) = p(y|\mathbf{x}, f) = \text{multinoulli}[f(\mathbf{x})]$$

(Probabilistic) Discriminative Classification

- Models $p(y|\mathbf{x})$ directly using an appropriate **discrete** distribution
- Parameters of the distribution $p(y|\mathbf{x})$ defined by a **linear/nonlinear function** of \mathbf{x}
 - A discriminative model only needs to learn these parameters!
- For binary classification, $p(y|\mathbf{x})$ would be a **Bernoulli**

$$p(y|\mathbf{x}) = p(y|\mathbf{x}, f) = \text{Bernoulli}[f(\mathbf{x})]$$

where $f()$ is a function (**defined by some params**) that maps \mathbf{x} to a probability $\mu = p(y = 1|\mathbf{x})$

- For $K > 2$ class multiclass classification, $p(y|\mathbf{x})$ would be a **multinoulli**

$$p(y|\mathbf{x}) = p(y|\mathbf{x}, f) = \text{multinoulli}[f(\mathbf{x})]$$

where $f()$ maps \mathbf{x} to a probability vector $\boldsymbol{\mu} = [\mu_1, \dots, \mu_K]$ s.t. $p(y = k|\mathbf{x}) = \mu_k$

(Probabilistic) Discriminative Classification

- Models $p(y|\mathbf{x})$ directly using an appropriate **discrete** distribution
- Parameters of the distribution $p(y|\mathbf{x})$ defined by a **linear/nonlinear function** of \mathbf{x}
 - A discriminative model only needs to learn these parameters!
- For binary classification, $p(y|\mathbf{x})$ would be a **Bernoulli**

$$p(y|\mathbf{x}) = p(y|\mathbf{x}, f) = \text{Bernoulli}[f(\mathbf{x})]$$

where $f()$ is a function (**defined by some params**) that maps \mathbf{x} to a probability $\mu = p(y = 1|\mathbf{x})$

- For $K > 2$ class multiclass classification, $p(y|\mathbf{x})$ would be a **multinoulli**

$$p(y|\mathbf{x}) = p(y|\mathbf{x}, f) = \text{multinoulli}[f(\mathbf{x})]$$

where $f()$ maps \mathbf{x} to a probability vector $\boldsymbol{\mu} = [\mu_1, \dots, \mu_K]$ s.t. $p(y = k|\mathbf{x}) = \mu_k$

- **Many choices of $f(\mathbf{x})$ possible.** Lead to different types of discriminative classification models

(Probabilistic) Discriminative Classification

- The general form of $p(y|\mathbf{x})$ in these models is given by

$$p(y|\mathbf{x}) = \text{Bernoulli}[f(\mathbf{x})] \quad \text{OR} \quad p(y|\mathbf{x}) = \text{multinoulli}[f(\mathbf{x})]$$

(Probabilistic) Discriminative Classification

- The general form of $p(y|\mathbf{x})$ in these models is given by

$$p(y|\mathbf{x}) = \text{Bernoulli}[f(\mathbf{x})] \quad \text{OR} \quad p(y|\mathbf{x}) = \text{multinoulli}[f(\mathbf{x})]$$

- The function $f(\mathbf{x})$ is **composed of two operations**

(Probabilistic) Discriminative Classification

- The general form of $p(y|\mathbf{x})$ in these models is given by

$$p(y|\mathbf{x}) = \text{Bernoulli}[f(\mathbf{x})] \quad \text{OR} \quad p(y|\mathbf{x}) = \text{multinoulli}[f(\mathbf{x})]$$

- The function $f(\mathbf{x})$ is composed of two operations
 - Compute a score for \mathbf{x} (e.g., using a linear model $\mathbf{w}^\top \mathbf{x}$ where \mathbf{w} is the param. to be estimated)

(Probabilistic) Discriminative Classification

- The general form of $p(y|\mathbf{x})$ in these models is given by

$$p(y|\mathbf{x}) = \text{Bernoulli}[f(\mathbf{x})] \quad \text{OR} \quad p(y|\mathbf{x}) = \text{multinoulli}[f(\mathbf{x})]$$

- The function $f(\mathbf{x})$ is composed of two operations
 - Compute a score for \mathbf{x} (e.g., using a linear model $\mathbf{w}^\top \mathbf{x}$ where \mathbf{w} is the param. to be estimated)
 - Turn this score into a probability of \mathbf{x} belonging to each class

(Probabilistic) Discriminative Classification

- The general form of $p(y|\mathbf{x})$ in these models is given by

$$p(y|\mathbf{x}) = \text{Bernoulli}[f(\mathbf{x})] \quad \text{OR} \quad p(y|\mathbf{x}) = \text{multinoulli}[f(\mathbf{x})]$$

- The function $f(\mathbf{x})$ is composed of two operations
 - Compute a **score** for \mathbf{x} (e.g., using a **linear model** $\mathbf{w}^\top \mathbf{x}$ where \mathbf{w} is the param. to be estimated)
 - Turn this score into a **probability** of \mathbf{x} belonging to each class
- Example: $f(\mathbf{x}) = \mu = \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})} \Rightarrow$ **Logistic Regression**; used for binary classification

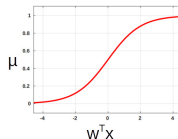
(Probabilistic) Discriminative Classification

- The general form of $p(y|\mathbf{x})$ in these models is given by

$$p(y|\mathbf{x}) = \text{Bernoulli}[f(\mathbf{x})] \quad \text{OR} \quad p(y|\mathbf{x}) = \text{multinoulli}[f(\mathbf{x})]$$

- The function $f(\mathbf{x})$ is **composed of two operations**
 - Compute a **score** for \mathbf{x} (e.g., using a **linear model** $\mathbf{w}^\top \mathbf{x}$ where \mathbf{w} is the param. to be estimated)
 - Turn this score into a **probability** of \mathbf{x} belonging to each class
- Example: $f(\mathbf{x}) = \mu = \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})} \Rightarrow$ **Logistic Regression**; used for binary classification

$$\mu = \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})}$$



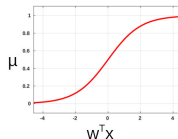
(Probabilistic) Discriminative Classification

- The general form of $p(y|\mathbf{x})$ in these models is given by

$$p(y|\mathbf{x}) = \text{Bernoulli}[f(\mathbf{x})] \quad \text{OR} \quad p(y|\mathbf{x}) = \text{multinoulli}[f(\mathbf{x})]$$

- The function $f(\mathbf{x})$ is **composed of two operations**
 - Compute a **score** for \mathbf{x} (e.g., using a **linear model** $\mathbf{w}^\top \mathbf{x}$ where \mathbf{w} is the param. to be estimated)
 - Turn this score into a **probability** of \mathbf{x} belonging to each class
- Example: $f(\mathbf{x}) = \mu = \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})} \Rightarrow$ **Logistic Regression**; used for binary classification

$$\mu = \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})}$$



- Another example: $f(\mathbf{x}) = \mu$ where $\mu_k = \frac{\exp(\mathbf{w}_k^\top \mathbf{x})}{\sum_{\ell=1}^K \exp(\mathbf{w}_\ell^\top \mathbf{x})} \Rightarrow$ **Softmax Regression**; used for multiclass

(Probabilistic) Discriminative Classification

- The general form of $p(y|\mathbf{x})$ in these models is given by

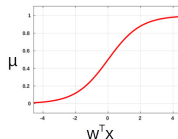
$$p(y|\mathbf{x}) = \text{Bernoulli}[f(\mathbf{x})] \quad \text{OR} \quad p(y|\mathbf{x}) = \text{multinoulli}[f(\mathbf{x})]$$

- The function $f(\mathbf{x})$ is **composed of two operations**

- Compute a **score** for \mathbf{x} (e.g., using a **linear model** $\mathbf{w}^\top \mathbf{x}$ where \mathbf{w} is the param. to be estimated)
- Turn this score into a **probability** of \mathbf{x} belonging to each class

- Example: $f(\mathbf{x}) = \mu = \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})} \Rightarrow$ **Logistic Regression**; used for binary classification

$$\mu = \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})}$$



- Another example: $f(\mathbf{x}) = \boldsymbol{\mu}$ where $\mu_k = \frac{\exp(\mathbf{w}_k^\top \mathbf{x})}{\sum_{\ell=1}^K \exp(\mathbf{w}_\ell^\top \mathbf{x})} \Rightarrow$ **Softmax Regression**; used for multiclass
- Can also use **nonlinear models** to compute the scores (e.g., deep NN or Gaussian Process)

Logistic Regression

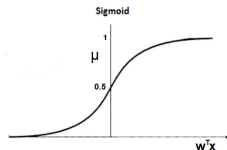
- Perhaps the simplest discriminative model for **linear binary classification**
- Defines $\mu = p(y = 1|\mathbf{x})$ using the **sigmoid function**

$$\mu = \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})} = \frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})}$$

Logistic Regression

- Perhaps the simplest discriminative model for **linear binary classification**
- Defines $\mu = p(y = 1|\mathbf{x})$ using the **sigmoid function**

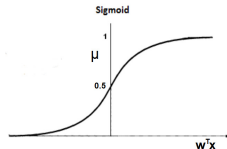
$$\mu = \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})} = \frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})}$$



Logistic Regression

- Perhaps the simplest discriminative model for linear binary classification
- Defines $\mu = p(y = 1|\mathbf{x})$ using the sigmoid function

$$\mu = \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})} = \frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})}$$

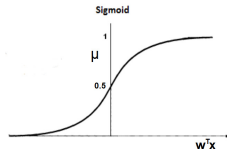


- Here $\mathbf{w}^\top \mathbf{x}$ is the score for input \mathbf{x} .

Logistic Regression

- Perhaps the simplest discriminative model for **linear binary classification**
- Defines $\mu = p(y = 1|\mathbf{x})$ using the **sigmoid function**

$$\mu = \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})} = \frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})}$$

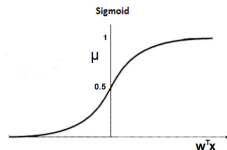


- Here $\mathbf{w}^\top \mathbf{x}$ is the score for input \mathbf{x} . The sigmoid turns it into a probability.

Logistic Regression

- Perhaps the simplest discriminative model for **linear binary classification**
- Defines $\mu = p(y = 1|\mathbf{x})$ using the **sigmoid function**

$$\mu = \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})} = \frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})}$$



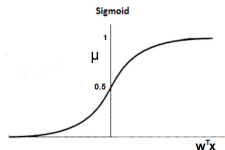
- Here $\mathbf{w}^\top \mathbf{x}$ is the score for input \mathbf{x} . The sigmoid turns it into a probability. Thus we have

$$\begin{aligned} p(y = 1|\mathbf{x}, \mathbf{w}) &= \mu = \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})} = \frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})} \\ p(y = 0|\mathbf{x}, \mathbf{w}) &= 1 - \mu = 1 - \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(\mathbf{w}^\top \mathbf{x})} \end{aligned}$$

Logistic Regression

- Perhaps the simplest discriminative model for linear binary classification
- Defines $\mu = p(y = 1|\mathbf{x})$ using the sigmoid function

$$\mu = \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})} = \frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})}$$



- Here $\mathbf{w}^\top \mathbf{x}$ is the score for input \mathbf{x} . The sigmoid turns it into a probability. Thus we have

$$\begin{aligned} p(y = 1|\mathbf{x}, \mathbf{w}) &= \mu = \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})} = \frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})} \\ p(y = 0|\mathbf{x}, \mathbf{w}) &= 1 - \mu = 1 - \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(\mathbf{w}^\top \mathbf{x})} \end{aligned}$$

- **Note:** If we assume $y \in \{-1, +1\}$ instead of $y \in \{0, 1\}$ then $p(y|\mathbf{x}, \mathbf{w}) = \frac{1}{1 + \exp(-y\mathbf{w}^\top \mathbf{x})}$

Logistic Regression: A Closer Look..

- At the decision boundary where both classes are equiprobable:

$$\begin{aligned}p(y = 1|\mathbf{x}, \mathbf{w}) &= p(y = 0|\mathbf{x}, \mathbf{w}) \\ \frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})} &= \frac{1}{1 + \exp(\mathbf{w}^\top \mathbf{x})} \\ \exp(\mathbf{w}^\top \mathbf{x}) &= 1 \\ \mathbf{w}^\top \mathbf{x} &= 0\end{aligned}$$

Logistic Regression: A Closer Look..

- At the decision boundary where both classes are equiprobable:

$$\begin{aligned}p(y = 1|\mathbf{x}, \mathbf{w}) &= p(y = 0|\mathbf{x}, \mathbf{w}) \\ \frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})} &= \frac{1}{1 + \exp(\mathbf{w}^\top \mathbf{x})} \\ \exp(\mathbf{w}^\top \mathbf{x}) &= 1 \\ \mathbf{w}^\top \mathbf{x} &= 0\end{aligned}$$

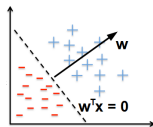
- Thus the decision boundary of LR is a **linear hyperplane**, just like Perceptron, SVM, etc.

Logistic Regression: A Closer Look..

- At the decision boundary where both classes are equiprobable:

$$\begin{aligned} p(y = 1 | \mathbf{x}, \mathbf{w}) &= p(y = 0 | \mathbf{x}, \mathbf{w}) \\ \frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})} &= \frac{1}{1 + \exp(\mathbf{w}^\top \mathbf{x})} \\ \exp(\mathbf{w}^\top \mathbf{x}) &= 1 \\ \mathbf{w}^\top \mathbf{x} &= 0 \end{aligned}$$

- Thus the decision boundary of LR is a **linear hyperplane**, just like Perceptron, SVM, etc.
- Therefore $y = 1$ if $\mathbf{w}^\top \mathbf{x} \geq 0$, otherwise $y = 0$

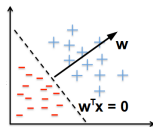


Logistic Regression: A Closer Look..

- At the decision boundary where both classes are equiprobable:

$$\begin{aligned}p(y = 1|\mathbf{x}, \mathbf{w}) &= p(y = 0|\mathbf{x}, \mathbf{w}) \\ \frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})} &= \frac{1}{1 + \exp(\mathbf{w}^\top \mathbf{x})} \\ \exp(\mathbf{w}^\top \mathbf{x}) &= 1 \\ \mathbf{w}^\top \mathbf{x} &= 0\end{aligned}$$

- Thus the decision boundary of LR is a **linear hyperplane**, just like Perceptron, SVM, etc.
- Therefore $y = 1$ if $\mathbf{w}^\top \mathbf{x} \geq 0$, otherwise $y = 0$



- High positive (negative) score $\mathbf{w}^\top \mathbf{x}$: High (low) probability of label 1

MLE for Logistic Regression

- Each label y_n is binary with probability $\mu_n = \frac{\exp(\mathbf{w}^\top \mathbf{x}_n)}{1 + \exp(\mathbf{w}^\top \mathbf{x}_n)}$. Since the likelihood is Bernoulli:

$$p(\mathbf{y}|\mathbf{X}, \mathbf{w}) = \prod_{n=1}^N p(y_n|\mathbf{x}_n, \mathbf{w}) = \prod_{n=1}^N \mu_n^{y_n} (1 - \mu_n)^{1-y_n}$$

MLE for Logistic Regression

- Each label y_n is binary with probability $\mu_n = \frac{\exp(\mathbf{w}^\top \mathbf{x}_n)}{1 + \exp(\mathbf{w}^\top \mathbf{x}_n)}$. Since the likelihood is Bernoulli:

$$p(\mathbf{y}|\mathbf{X}, \mathbf{w}) = \prod_{n=1}^N p(y_n|\mathbf{x}_n, \mathbf{w}) = \prod_{n=1}^N \mu_n^{y_n} (1 - \mu_n)^{1-y_n}$$

- Negative log-likelihood

$$\text{NLL}(\mathbf{w}) = -\log p(\mathbf{y}|\mathbf{X}, \mathbf{w}) = -\sum_{n=1}^N (y_n \log \mu_n + (1 - y_n) \log(1 - \mu_n))$$

MLE for Logistic Regression

- Each label y_n is binary with probability $\mu_n = \frac{\exp(\mathbf{w}^\top \mathbf{x}_n)}{1 + \exp(\mathbf{w}^\top \mathbf{x}_n)}$. Since the likelihood is Bernoulli:

$$p(\mathbf{y}|\mathbf{X}, \mathbf{w}) = \prod_{n=1}^N p(y_n|\mathbf{x}_n, \mathbf{w}) = \prod_{n=1}^N \mu_n^{y_n} (1 - \mu_n)^{1-y_n}$$

- Negative log-likelihood

$$\text{NLL}(\mathbf{w}) = -\log p(\mathbf{y}|\mathbf{X}, \mathbf{w}) = -\sum_{n=1}^N (y_n \log \mu_n + (1 - y_n) \log(1 - \mu_n))$$

- Plugging in $\mu_n = \frac{\exp(\mathbf{w}^\top \mathbf{x}_n)}{1 + \exp(\mathbf{w}^\top \mathbf{x}_n)}$ and chugging, we get (verify yourself)

$$\text{NLL}(\mathbf{w}) = -\sum_{n=1}^N (y_n \mathbf{w}^\top \mathbf{x}_n - \log(1 + \exp(\mathbf{w}^\top \mathbf{x}_n)))$$

MLE for Logistic Regression

- Each label y_n is binary with probability $\mu_n = \frac{\exp(\mathbf{w}^\top \mathbf{x}_n)}{1 + \exp(\mathbf{w}^\top \mathbf{x}_n)}$. Since the likelihood is Bernoulli:

$$p(\mathbf{y}|\mathbf{X}, \mathbf{w}) = \prod_{n=1}^N p(y_n|\mathbf{x}_n, \mathbf{w}) = \prod_{n=1}^N \mu_n^{y_n} (1 - \mu_n)^{1-y_n}$$

- Negative log-likelihood

$$\text{NLL}(\mathbf{w}) = -\log p(\mathbf{y}|\mathbf{X}, \mathbf{w}) = -\sum_{n=1}^N (y_n \log \mu_n + (1 - y_n) \log(1 - \mu_n))$$

- Plugging in $\mu_n = \frac{\exp(\mathbf{w}^\top \mathbf{x}_n)}{1 + \exp(\mathbf{w}^\top \mathbf{x}_n)}$ and chugging, we get (verify yourself)

$$\text{NLL}(\mathbf{w}) = -\sum_{n=1}^N (y_n \mathbf{w}^\top \mathbf{x}_n - \log(1 + \exp(\mathbf{w}^\top \mathbf{x}_n)))$$

- MLE solution: $\mathbf{w}_{MLE} = \arg \min_{\mathbf{w}} \text{NLL}(\mathbf{w})$

MLE for Logistic Regression

- Each label y_n is binary with probability $\mu_n = \frac{\exp(\mathbf{w}^\top \mathbf{x}_n)}{1 + \exp(\mathbf{w}^\top \mathbf{x}_n)}$. Since the likelihood is Bernoulli:

$$p(\mathbf{y}|\mathbf{X}, \mathbf{w}) = \prod_{n=1}^N p(y_n|\mathbf{x}_n, \mathbf{w}) = \prod_{n=1}^N \mu_n^{y_n} (1 - \mu_n)^{1-y_n}$$

- Negative log-likelihood

$$\text{NLL}(\mathbf{w}) = -\log p(\mathbf{y}|\mathbf{X}, \mathbf{w}) = -\sum_{n=1}^N (y_n \log \mu_n + (1 - y_n) \log(1 - \mu_n))$$

- Plugging in $\mu_n = \frac{\exp(\mathbf{w}^\top \mathbf{x}_n)}{1 + \exp(\mathbf{w}^\top \mathbf{x}_n)}$ and chugging, we get (verify yourself)

$$\text{NLL}(\mathbf{w}) = -\sum_{n=1}^N (y_n \mathbf{w}^\top \mathbf{x}_n - \log(1 + \exp(\mathbf{w}^\top \mathbf{x}_n)))$$

- MLE solution: $\mathbf{w}_{MLE} = \arg \min_{\mathbf{w}} \text{NLL}(\mathbf{w})$
- Important note:** $\text{NLL}(\mathbf{w})$ is convex in \mathbf{w} , so global minima can be found

MLE for Logistic Regression

- We have $\text{NLL}(\mathbf{w}) = -\sum_{n=1}^N (y_n \mathbf{w}^\top \mathbf{x}_n - \log(1 + \exp(\mathbf{w}^\top \mathbf{x}_n)))$

MLE for Logistic Regression

- We have $\text{NLL}(\mathbf{w}) = -\sum_{n=1}^N (y_n \mathbf{w}^\top \mathbf{x}_n - \log(1 + \exp(\mathbf{w}^\top \mathbf{x}_n)))$
- Taking the derivative of $\text{NLL}(\mathbf{w})$ w.r.t. \mathbf{w}

$$\begin{aligned}\frac{\partial \text{NLL}(\mathbf{w})}{\partial \mathbf{w}} &= \frac{\partial}{\partial \mathbf{w}} \left[-\sum_{n=1}^N (y_n \mathbf{w}^\top \mathbf{x}_n - \log(1 + \exp(\mathbf{w}^\top \mathbf{x}_n))) \right] \\ &= -\sum_{n=1}^N \left(y_n \mathbf{x}_n - \frac{\exp(\mathbf{w}^\top \mathbf{x}_n)}{(1 + \exp(\mathbf{w}^\top \mathbf{x}_n))} \mathbf{x}_n \right)\end{aligned}$$

MLE for Logistic Regression

- We have $\text{NLL}(\mathbf{w}) = -\sum_{n=1}^N (y_n \mathbf{w}^\top \mathbf{x}_n - \log(1 + \exp(\mathbf{w}^\top \mathbf{x}_n)))$
- Taking the derivative of $\text{NLL}(\mathbf{w})$ w.r.t. \mathbf{w}

$$\begin{aligned}\frac{\partial \text{NLL}(\mathbf{w})}{\partial \mathbf{w}} &= \frac{\partial}{\partial \mathbf{w}} \left[-\sum_{n=1}^N (y_n \mathbf{w}^\top \mathbf{x}_n - \log(1 + \exp(\mathbf{w}^\top \mathbf{x}_n))) \right] \\ &= -\sum_{n=1}^N \left(y_n \mathbf{x}_n - \frac{\exp(\mathbf{w}^\top \mathbf{x}_n)}{(1 + \exp(\mathbf{w}^\top \mathbf{x}_n))} \mathbf{x}_n \right)\end{aligned}$$

- Can't get a closed form estimate for \mathbf{w} by setting the derivative to zero

MLE for Logistic Regression

- We have $\text{NLL}(\mathbf{w}) = -\sum_{n=1}^N (y_n \mathbf{w}^\top \mathbf{x}_n - \log(1 + \exp(\mathbf{w}^\top \mathbf{x}_n)))$
- Taking the derivative of $\text{NLL}(\mathbf{w})$ w.r.t. \mathbf{w}

$$\begin{aligned}\frac{\partial \text{NLL}(\mathbf{w})}{\partial \mathbf{w}} &= \frac{\partial}{\partial \mathbf{w}} \left[-\sum_{n=1}^N (y_n \mathbf{w}^\top \mathbf{x}_n - \log(1 + \exp(\mathbf{w}^\top \mathbf{x}_n))) \right] \\ &= -\sum_{n=1}^N \left(y_n \mathbf{x}_n - \frac{\exp(\mathbf{w}^\top \mathbf{x}_n)}{(1 + \exp(\mathbf{w}^\top \mathbf{x}_n))} \mathbf{x}_n \right)\end{aligned}$$

- Can't get a closed form estimate for \mathbf{w} by setting the derivative to zero
- One solution: Iterative minimization via gradient descent $\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \mathbf{g}_t$. The gradient is:

$$\mathbf{g} = \frac{\partial \text{NLL}(\mathbf{w})}{\partial \mathbf{w}} = -\sum_{n=1}^N (y_n - \mu_n) \mathbf{x}_n = \mathbf{X}^\top (\boldsymbol{\mu} - \mathbf{y})$$

MLE for Logistic Regression

- We have $\text{NLL}(\mathbf{w}) = -\sum_{n=1}^N (y_n \mathbf{w}^\top \mathbf{x}_n - \log(1 + \exp(\mathbf{w}^\top \mathbf{x}_n)))$
- Taking the derivative of $\text{NLL}(\mathbf{w})$ w.r.t. \mathbf{w}

$$\begin{aligned}\frac{\partial \text{NLL}(\mathbf{w})}{\partial \mathbf{w}} &= \frac{\partial}{\partial \mathbf{w}} \left[-\sum_{n=1}^N (y_n \mathbf{w}^\top \mathbf{x}_n - \log(1 + \exp(\mathbf{w}^\top \mathbf{x}_n))) \right] \\ &= -\sum_{n=1}^N \left(y_n \mathbf{x}_n - \frac{\exp(\mathbf{w}^\top \mathbf{x}_n)}{(1 + \exp(\mathbf{w}^\top \mathbf{x}_n))} \mathbf{x}_n \right)\end{aligned}$$

- Can't get a closed form estimate for \mathbf{w} by setting the derivative to zero
- One solution: Iterative minimization via gradient descent $\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \mathbf{g}_t$. The gradient is:

$$\mathbf{g} = \frac{\partial \text{NLL}(\mathbf{w})}{\partial \mathbf{w}} = -\sum_{n=1}^N (y_n - \mu_n) \mathbf{x}_n = \mathbf{X}^\top (\boldsymbol{\mu} - \mathbf{y})$$

- A large error on $\mathbf{x}_n \Rightarrow (y_n - \mu_n)$ will be large \Rightarrow large contribution of \mathbf{x}_n to the gradient

MLE for Logistic Regression

- We have $\text{NLL}(\mathbf{w}) = -\sum_{n=1}^N (y_n \mathbf{w}^\top \mathbf{x}_n - \log(1 + \exp(\mathbf{w}^\top \mathbf{x}_n)))$
- Taking the derivative of $\text{NLL}(\mathbf{w})$ w.r.t. \mathbf{w}

$$\begin{aligned}\frac{\partial \text{NLL}(\mathbf{w})}{\partial \mathbf{w}} &= \frac{\partial}{\partial \mathbf{w}} \left[-\sum_{n=1}^N (y_n \mathbf{w}^\top \mathbf{x}_n - \log(1 + \exp(\mathbf{w}^\top \mathbf{x}_n))) \right] \\ &= -\sum_{n=1}^N \left(y_n \mathbf{x}_n - \frac{\exp(\mathbf{w}^\top \mathbf{x}_n)}{(1 + \exp(\mathbf{w}^\top \mathbf{x}_n))} \mathbf{x}_n \right)\end{aligned}$$

- Can't get a closed form estimate for \mathbf{w} by setting the derivative to zero
- One solution: Iterative minimization via gradient descent $\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \mathbf{g}_t$. The gradient is:

$$\mathbf{g} = \frac{\partial \text{NLL}(\mathbf{w})}{\partial \mathbf{w}} = -\sum_{n=1}^N (y_n - \mu_n) \mathbf{x}_n = \mathbf{X}^\top (\boldsymbol{\mu} - \mathbf{y})$$

- A large error on $\mathbf{x}_n \Rightarrow (y_n - \mu_n)$ will be large \Rightarrow large contribution of \mathbf{x}_n to the gradient
- More sophisticated methods (e.g., Newton's method) can also be used (I'll provide a note)

MAP Estimation for Logistic Regression

- MLE estimate of \mathbf{w} can lead to overfitting. Solution: use a prior $p(\mathbf{w})$ on \mathbf{w}
- Just like the linear regression case, let's put a Gaussian prior on \mathbf{w}

$$p(\mathbf{w}) = \mathcal{N}(0, \lambda^{-1} \mathbf{I}_D) \propto \exp\left(-\frac{\lambda}{2} \mathbf{w}^\top \mathbf{w}\right)$$

MAP Estimation for Logistic Regression

- MLE estimate of \mathbf{w} can lead to overfitting. Solution: use a prior $p(\mathbf{w})$ on \mathbf{w}
- Just like the linear regression case, let's put a Gaussian prior on \mathbf{w}

$$p(\mathbf{w}) = \mathcal{N}(0, \lambda^{-1} \mathbf{I}_D) \propto \exp\left(-\frac{\lambda}{2} \mathbf{w}^\top \mathbf{w}\right)$$

- MAP objective (log of posterior) = MLE objective + $\log p(\mathbf{w})$
- Leads to the objective (negative of log posterior, ignoring constants):

$$\text{NLL}(\mathbf{w}) + \frac{\lambda}{2} \mathbf{w}^\top \mathbf{w}$$

MAP Estimation for Logistic Regression

- MLE estimate of \mathbf{w} can lead to overfitting. Solution: use a prior $p(\mathbf{w})$ on \mathbf{w}
- Just like the linear regression case, let's put a Gaussian prior on \mathbf{w}

$$p(\mathbf{w}) = \mathcal{N}(0, \lambda^{-1} \mathbf{I}_D) \propto \exp\left(-\frac{\lambda}{2} \mathbf{w}^\top \mathbf{w}\right)$$

- MAP objective (log of posterior) = MLE objective + $\log p(\mathbf{w})$
- Leads to the objective (negative of log posterior, ignoring constants):

$$\text{NLL}(\mathbf{w}) + \frac{\lambda}{2} \mathbf{w}^\top \mathbf{w}$$

- Estimation of \mathbf{w}_{MAP} proceeds the same way as MLE using gradient based methods, with

$$\text{Gradient: } \mathbf{g} = \mathbf{X}^\top (\boldsymbol{\mu} - \mathbf{y}) + \lambda \mathbf{w}$$

Fully Bayesian Estimation for Logistic Regression

- Doing fully Bayesian inference would require computing the posterior

$$p(\mathbf{w}|\mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w})}{\int p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w})d\mathbf{w}}$$

Fully Bayesian Estimation for Logistic Regression

- Doing fully Bayesian inference would require computing the posterior

$$p(\mathbf{w}|\mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w})}{\int p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w})d\mathbf{w}} = \frac{\prod_{n=1}^N p(y_n|\mathbf{x}_n, \mathbf{w})p(\mathbf{w})}{\int \prod_{n=1}^N p(y_n|\mathbf{x}_n, \mathbf{w})p(\mathbf{w})d\mathbf{w}}$$

- **Intractable**. Reason: likelihood (logistic-Bernoulli) and prior (Gaussian) here are **not conjugate**

Fully Bayesian Estimation for Logistic Regression

- Doing fully Bayesian inference would require computing the posterior

$$p(\mathbf{w}|\mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w})}{\int p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w})d\mathbf{w}} = \frac{\prod_{n=1}^N p(y_n|\mathbf{x}_n, \mathbf{w})p(\mathbf{w})}{\int \prod_{n=1}^N p(y_n|\mathbf{x}_n, \mathbf{w})p(\mathbf{w})d\mathbf{w}}$$

- **Intractable**. Reason: likelihood (logistic-Bernoulli) and prior (Gaussian) here are **not conjugate**
- Need to do *approximate inference* in this case

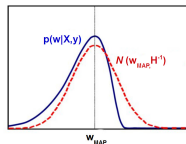
Fully Bayesian Estimation for Logistic Regression

- Doing fully Bayesian inference would require computing the posterior

$$p(\mathbf{w}|\mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w})}{\int p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w})d\mathbf{w}} = \frac{\prod_{n=1}^N p(y_n|\mathbf{x}_n, \mathbf{w})p(\mathbf{w})}{\int \prod_{n=1}^N p(y_n|\mathbf{x}_n, \mathbf{w})p(\mathbf{w})d\mathbf{w}}$$

- **Intractable**. Reason: likelihood (logistic-Bernoulli) and prior (Gaussian) here are **not conjugate**
- Need to do **approximate inference** in this case
- A crude approximation: **Laplace approximation**: Approximate a posterior by a **Gaussian** with **mean** = \mathbf{w}_{MAP} and **covariance** = **inverse hessian** (hessian = second derivative of $\log p(\mathbf{w}|\mathbf{X}, \mathbf{y})$)

$$p(\mathbf{w}|\mathbf{X}, \mathbf{y}) = \mathcal{N}(\mathbf{w}_{MAP}, \mathbf{H}^{-1})$$



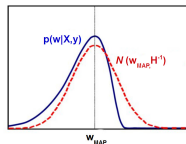
Fully Bayesian Estimation for Logistic Regression

- Doing fully Bayesian inference would require computing the posterior

$$p(\mathbf{w}|\mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w})}{\int p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w})d\mathbf{w}} = \frac{\prod_{n=1}^N p(y_n|\mathbf{x}_n, \mathbf{w})p(\mathbf{w})}{\int \prod_{n=1}^N p(y_n|\mathbf{x}_n, \mathbf{w})p(\mathbf{w})d\mathbf{w}}$$

- **Intractable**. Reason: likelihood (logistic-Bernoulli) and prior (Gaussian) here are **not conjugate**
- Need to do **approximate inference** in this case
- A crude approximation: **Laplace approximation**: Approximate a posterior by a **Gaussian** with **mean** = \mathbf{w}_{MAP} and **covariance** = **inverse hessian** (hessian = second derivative of $\log p(\mathbf{w}|\mathbf{X}, \mathbf{y})$)

$$p(\mathbf{w}|\mathbf{X}, \mathbf{y}) = \mathcal{N}(\mathbf{w}_{MAP}, \mathbf{H}^{-1})$$



- Will look at Laplace and other approximate inference methods later during the semester

Logistic Regression: Predictive Distributions

- When using MLE, the predictive distribution will be

$$p(y_* = 1 | \mathbf{x}_*, \mathbf{w}_{MLE}) = \sigma(\mathbf{w}_{MLE}^\top \mathbf{x}_*)$$

Logistic Regression: Predictive Distributions

- When using MLE, the predictive distribution will be

$$\begin{aligned} p(y_* = 1 | \mathbf{x}_*, \mathbf{w}_{MLE}) &= \sigma(\mathbf{w}_{MLE}^\top \mathbf{x}_*) \\ p(y_* | \mathbf{x}_*, \mathbf{w}_{MLE}) &= \text{Bernoulli}(\sigma(\mathbf{w}_{MLE}^\top \mathbf{x}_*)) \end{aligned}$$

Logistic Regression: Predictive Distributions

- When using MLE, the predictive distribution will be

$$\begin{aligned}p(y_* = 1 | \mathbf{x}_*, \mathbf{w}_{MLE}) &= \sigma(\mathbf{w}_{MLE}^\top \mathbf{x}_*) \\p(y_* | \mathbf{x}_*, \mathbf{w}_{MLE}) &= \text{Bernoulli}(\sigma(\mathbf{w}_{MLE}^\top \mathbf{x}_*))\end{aligned}$$

- When using MAP, the predictive distribution will be

$$p(y_* | \mathbf{x}_*, \mathbf{w}_{MAP}) = \text{Bernoulli}(\sigma(\mathbf{w}_{MAP}^\top \mathbf{x}_*))$$

Logistic Regression: Predictive Distributions

- When using MLE, the predictive distribution will be

$$\begin{aligned}p(y_* = 1 | \mathbf{x}_*, \mathbf{w}_{MLE}) &= \sigma(\mathbf{w}_{MLE}^\top \mathbf{x}_*) \\p(y_* | \mathbf{x}_*, \mathbf{w}_{MLE}) &= \text{Bernoulli}(\sigma(\mathbf{w}_{MLE}^\top \mathbf{x}_*))\end{aligned}$$

- When using MAP, the predictive distribution will be

$$\begin{aligned}p(y_* | \mathbf{x}_*, \mathbf{w}_{MAP}) &= \text{Bernoulli}(\sigma(\mathbf{w}_{MAP}^\top \mathbf{x}_*)) \\p(y_* = 1 | \mathbf{x}_*, \mathbf{w}_{MAP}) &= \sigma(\mathbf{w}_{MAP}^\top \mathbf{x}_*)\end{aligned}$$

Logistic Regression: Predictive Distributions

- When using MLE, the predictive distribution will be

$$\begin{aligned}p(y_* = 1 | \mathbf{x}_*, \mathbf{w}_{MLE}) &= \sigma(\mathbf{w}_{MLE}^\top \mathbf{x}_*) \\p(y_* | \mathbf{x}_*, \mathbf{w}_{MLE}) &= \text{Bernoulli}(\sigma(\mathbf{w}_{MLE}^\top \mathbf{x}_*))\end{aligned}$$

- When using MAP, the predictive distribution will be

$$\begin{aligned}p(y_* | \mathbf{x}_*, \mathbf{w}_{MAP}) &= \text{Bernoulli}(\sigma(\mathbf{w}_{MAP}^\top \mathbf{x}_*)) \\p(y_* = 1 | \mathbf{x}_*, \mathbf{w}_{MAP}) &= \sigma(\mathbf{w}_{MAP}^\top \mathbf{x}_*)\end{aligned}$$

- When using Bayesian inference, the **posterior predictive distribution**, based on posterior averaging

$$p(y_* = 1 | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \int p(y_* = 1 | \mathbf{x}_*, \mathbf{w}) p(\mathbf{w} | \mathbf{X}, \mathbf{y}) d\mathbf{w}$$

Logistic Regression: Predictive Distributions

- When using MLE, the predictive distribution will be

$$\begin{aligned}p(y_* = 1 | \mathbf{x}_*, \mathbf{w}_{MLE}) &= \sigma(\mathbf{w}_{MLE}^\top \mathbf{x}_*) \\p(y_* | \mathbf{x}_*, \mathbf{w}_{MLE}) &= \text{Bernoulli}(\sigma(\mathbf{w}_{MLE}^\top \mathbf{x}_*))\end{aligned}$$

- When using MAP, the predictive distribution will be

$$\begin{aligned}p(y_* | \mathbf{x}_*, \mathbf{w}_{MAP}) &= \text{Bernoulli}(\sigma(\mathbf{w}_{MAP}^\top \mathbf{x}_*)) \\p(y_* = 1 | \mathbf{x}_*, \mathbf{w}_{MAP}) &= \sigma(\mathbf{w}_{MAP}^\top \mathbf{x}_*)\end{aligned}$$

- When using Bayesian inference, the **posterior predictive distribution**, based on posterior averaging

$$p(y_* = 1 | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \int p(y_* = 1 | \mathbf{x}_*, \mathbf{w}) p(\mathbf{w} | \mathbf{X}, \mathbf{y}) d\mathbf{w} = \int \sigma(\mathbf{w}^\top \mathbf{x}_*) p(\mathbf{w} | \mathbf{X}, \mathbf{y}) d\mathbf{w}$$

Logistic Regression: Predictive Distributions

- When using MLE, the predictive distribution will be

$$\begin{aligned}p(y_* = 1 | \mathbf{x}_*, \mathbf{w}_{MLE}) &= \sigma(\mathbf{w}_{MLE}^\top \mathbf{x}_*) \\p(y_* | \mathbf{x}_*, \mathbf{w}_{MLE}) &= \text{Bernoulli}(\sigma(\mathbf{w}_{MLE}^\top \mathbf{x}_*))\end{aligned}$$

- When using MAP, the predictive distribution will be

$$\begin{aligned}p(y_* | \mathbf{x}_*, \mathbf{w}_{MAP}) &= \text{Bernoulli}(\sigma(\mathbf{w}_{MAP}^\top \mathbf{x}_*)) \\p(y_* = 1 | \mathbf{x}_*, \mathbf{w}_{MAP}) &= \sigma(\mathbf{w}_{MAP}^\top \mathbf{x}_*)\end{aligned}$$

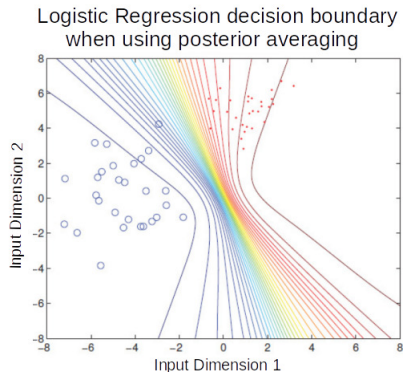
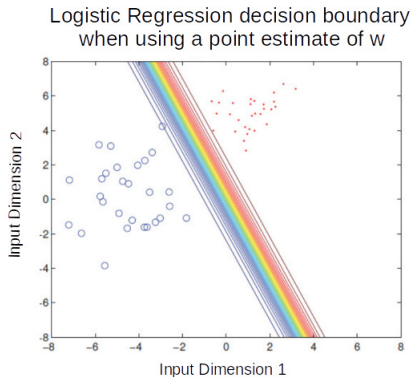
- When using Bayesian inference, the **posterior predictive distribution**, based on posterior averaging

$$p(y_* = 1 | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \int p(y_* = 1 | \mathbf{x}_*, \mathbf{w}) p(\mathbf{w} | \mathbf{X}, \mathbf{y}) d\mathbf{w} = \int \sigma(\mathbf{w}^\top \mathbf{x}_*) p(\mathbf{w} | \mathbf{X}, \mathbf{y}) d\mathbf{w}$$

- **Note:** Unlike the linear regression case, for logistic regression (**and for non-conjugate models in general**), posterior averaging can be intractable (and may require approximations)

Logistic Regression: Plug-in Prediction vs Bayesian Averaging

- (Left) Predictive distribution when using a point estimate uses only a single linear hyperplane \mathbf{w}
- (Right) Posterior predictive distribution **averages over many linear hyperplanes** \mathbf{w}



Multiclass Logistic (a.k.a. Softmax) Regression

- Also called multinoulli/multinomial regression: Basically, logistic regression for $K > 2$ classes

Multiclass Logistic (a.k.a. Softmax) Regression

- Also called multinoulli/multinomial regression: Basically, logistic regression for $K > 2$ classes
- In this case, $y_n \in \{1, 2, \dots, K\}$ and label probabilities are defined as

$$p(y_n = k | \mathbf{x}_n, \mathbf{W}) = \frac{\exp(\mathbf{w}_k^\top \mathbf{x}_n)}{\sum_{\ell=1}^K \exp(\mathbf{w}_\ell^\top \mathbf{x}_n)} = \mu_{nk}$$

Multiclass Logistic (a.k.a. Softmax) Regression

- Also called multinoulli/multinomial regression: Basically, logistic regression for $K > 2$ classes
- In this case, $y_n \in \{1, 2, \dots, K\}$ and label probabilities are defined as

$$p(y_n = k | \mathbf{x}_n, \mathbf{W}) = \frac{\exp(\mathbf{w}_k^\top \mathbf{x}_n)}{\sum_{\ell=1}^K \exp(\mathbf{w}_\ell^\top \mathbf{x}_n)} = \mu_{nk} \quad \text{and} \quad \sum_{\ell=1}^K \mu_{n\ell} = 1$$

Multiclass Logistic (a.k.a. Softmax) Regression

- Also called multinoulli/multinomial regression: Basically, logistic regression for $K > 2$ classes
- In this case, $y_n \in \{1, 2, \dots, K\}$ and label probabilities are defined as

$$p(y_n = k | \mathbf{x}_n, \mathbf{W}) = \frac{\exp(\mathbf{w}_k^\top \mathbf{x}_n)}{\sum_{\ell=1}^K \exp(\mathbf{w}_\ell^\top \mathbf{x}_n)} = \mu_{nk} \quad \text{and} \quad \sum_{\ell=1}^K \mu_{n\ell} = 1$$

- $\mathbf{W} = [\mathbf{w}_1 \ \mathbf{w}_2 \ \dots \ \mathbf{w}_K]$ is $D \times K$ weight matrix.

Multiclass Logistic (a.k.a. Softmax) Regression

- Also called multinoulli/multinomial regression: Basically, logistic regression for $K > 2$ classes
- In this case, $y_n \in \{1, 2, \dots, K\}$ and label probabilities are defined as

$$p(y_n = k | \mathbf{x}_n, \mathbf{W}) = \frac{\exp(\mathbf{w}_k^\top \mathbf{x}_n)}{\sum_{\ell=1}^K \exp(\mathbf{w}_\ell^\top \mathbf{x}_n)} = \mu_{nk} \quad \text{and} \quad \sum_{\ell=1}^K \mu_{n\ell} = 1$$

- $\mathbf{W} = [\mathbf{w}_1 \ \mathbf{w}_2 \ \dots \ \mathbf{w}_K]$ is $D \times K$ **weight matrix**. $\mathbf{w}_1 = \mathbf{0}_{D \times 1}$ (assumed for identifiability)

Multiclass Logistic (a.k.a. Softmax) Regression

- Also called multinoulli/multinomial regression: Basically, logistic regression for $K > 2$ classes
- In this case, $y_n \in \{1, 2, \dots, K\}$ and label probabilities are defined as

$$p(y_n = k | \mathbf{x}_n, \mathbf{W}) = \frac{\exp(\mathbf{w}_k^\top \mathbf{x}_n)}{\sum_{\ell=1}^K \exp(\mathbf{w}_\ell^\top \mathbf{x}_n)} = \mu_{nk} \quad \text{and} \quad \sum_{\ell=1}^K \mu_{n\ell} = 1$$

- $\mathbf{W} = [\mathbf{w}_1 \ \mathbf{w}_2 \ \dots \ \mathbf{w}_K]$ is $D \times K$ **weight matrix**. $\mathbf{w}_1 = \mathbf{0}_{D \times 1}$ (assumed for identifiability)
- Called **softmax** because.. ?

Multiclass Logistic (a.k.a. Softmax) Regression

- Also called multinoulli/multinomial regression: Basically, logistic regression for $K > 2$ classes
- In this case, $y_n \in \{1, 2, \dots, K\}$ and label probabilities are defined as

$$p(y_n = k | \mathbf{x}_n, \mathbf{W}) = \frac{\exp(\mathbf{w}_k^\top \mathbf{x}_n)}{\sum_{\ell=1}^K \exp(\mathbf{w}_\ell^\top \mathbf{x}_n)} = \mu_{nk} \quad \text{and} \quad \sum_{\ell=1}^K \mu_{n\ell} = 1$$

- $\mathbf{W} = [\mathbf{w}_1 \ \mathbf{w}_2 \ \dots \ \mathbf{w}_K]$ is $D \times K$ weight matrix. $\mathbf{w}_1 = \mathbf{0}_{D \times 1}$ (assumed for identifiability)
- Called softmax because.. ?
- Each likelihood $p(y_n | \mathbf{x}_n, \mathbf{W})$ is a multinoulli. Therefore

$$p(\mathbf{y} | \mathbf{X}, \mathbf{W}) = \prod_{n=1}^N \prod_{\ell=1}^K \mu_{n\ell}^{y_{n\ell}}$$

Multiclass Logistic (a.k.a. Softmax) Regression

- Also called multinoulli/multinomial regression: Basically, logistic regression for $K > 2$ classes
- In this case, $y_n \in \{1, 2, \dots, K\}$ and label probabilities are defined as

$$p(y_n = k | \mathbf{x}_n, \mathbf{W}) = \frac{\exp(\mathbf{w}_k^\top \mathbf{x}_n)}{\sum_{\ell=1}^K \exp(\mathbf{w}_\ell^\top \mathbf{x}_n)} = \mu_{nk} \quad \text{and} \quad \sum_{\ell=1}^K \mu_{n\ell} = 1$$

- $\mathbf{W} = [\mathbf{w}_1 \ \mathbf{w}_2 \ \dots \ \mathbf{w}_K]$ is $D \times K$ **weight matrix**. $\mathbf{w}_1 = \mathbf{0}_{D \times 1}$ (assumed for identifiability)
- Called **softmax** because.. ?
- Each likelihood $p(y_n | \mathbf{x}_n, \mathbf{W})$ is a **multinoulli**. Therefore

$$p(\mathbf{y} | \mathbf{X}, \mathbf{W}) = \prod_{n=1}^N \prod_{\ell=1}^K \mu_{n\ell}^{y_{n\ell}}$$

where $y_{n\ell} = 1$ if true class of example n is ℓ and $y_{n\ell'} = 0$ for all other $\ell' \neq \ell$

Multiclass Logistic (a.k.a. Softmax) Regression

- Also called multinoulli/multinomial regression: Basically, logistic regression for $K > 2$ classes
- In this case, $y_n \in \{1, 2, \dots, K\}$ and label probabilities are defined as

$$p(y_n = k | \mathbf{x}_n, \mathbf{W}) = \frac{\exp(\mathbf{w}_k^\top \mathbf{x}_n)}{\sum_{\ell=1}^K \exp(\mathbf{w}_\ell^\top \mathbf{x}_n)} = \mu_{nk} \quad \text{and} \quad \sum_{\ell=1}^K \mu_{n\ell} = 1$$

- $\mathbf{W} = [\mathbf{w}_1 \ \mathbf{w}_2 \ \dots \ \mathbf{w}_K]$ is $D \times K$ **weight matrix**. $\mathbf{w}_1 = \mathbf{0}_{D \times 1}$ (assumed for identifiability)
- Called **softmax** because.. ?
- Each likelihood $p(y_n | \mathbf{x}_n, \mathbf{W})$ is a **multinoulli**. Therefore

$$p(\mathbf{y} | \mathbf{X}, \mathbf{W}) = \prod_{n=1}^N \prod_{\ell=1}^K \mu_{n\ell}^{y_{n\ell}}$$

where $y_{n\ell} = 1$ if true class of example n is ℓ and $y_{n\ell'} = 0$ for all other $\ell' \neq \ell$

- Can do MLE/MAP/fully Bayesian estimation for \mathbf{W} similar to the logistic regression model

Multiclass Logistic (a.k.a. Softmax) Regression

- Also called multinoulli/multinomial regression: Basically, logistic regression for $K > 2$ classes
- In this case, $y_n \in \{1, 2, \dots, K\}$ and label probabilities are defined as

$$p(y_n = k | \mathbf{x}_n, \mathbf{W}) = \frac{\exp(\mathbf{w}_k^\top \mathbf{x}_n)}{\sum_{\ell=1}^K \exp(\mathbf{w}_\ell^\top \mathbf{x}_n)} = \mu_{nk} \quad \text{and} \quad \sum_{\ell=1}^K \mu_{n\ell} = 1$$

- $\mathbf{W} = [\mathbf{w}_1 \ \mathbf{w}_2 \ \dots \ \mathbf{w}_K]$ is $D \times K$ **weight matrix**. $\mathbf{w}_1 = \mathbf{0}_{D \times 1}$ (assumed for identifiability)
- Called **softmax** because.. ?
- Each likelihood $p(y_n | \mathbf{x}_n, \mathbf{W})$ is a **multinoulli**. Therefore

$$p(\mathbf{y} | \mathbf{X}, \mathbf{W}) = \prod_{n=1}^N \prod_{\ell=1}^K \mu_{n\ell}^{y_{n\ell}}$$

where $y_{n\ell} = 1$ if true class of example n is ℓ and $y_{n\ell'} = 0$ for all other $\ell' \neq \ell$

- Can do MLE/MAP/fully Bayesian estimation for \mathbf{W} similar to the logistic regression model
- Doing MLE is like minimizing the cross-entropy loss: $-\sum_{n=1}^N \sum_{\ell=1}^K y_{n\ell} \log \mu_{n\ell}$

Some Thoughts on Logistic and Softmax Classification

- Both logistic and softmax classification are discriminative model for linear classification

Some Thoughts on Logistic and Softmax Classification

- Both logistic and softmax classification are discriminative model for linear classification
- These are special cases of the more general form of discriminative models for classification

$$p(y|\mathbf{x}) = \text{Bernoulli}(f(\mathbf{x})) \quad \text{OR} \quad p(y|\mathbf{x}) = \text{multinoulli}(f(\mathbf{x}))$$

Some Thoughts on Logistic and Softmax Classification

- Both logistic and softmax classification are discriminative model for linear classification
- These are special cases of the more general form of discriminative models for classification

$$p(y|\mathbf{x}) = \text{Bernoulli}(f(\mathbf{x})) \quad \text{OR} \quad p(y|\mathbf{x}) = \text{multinoulli}(f(\mathbf{x}))$$

- Logistic and softmax also belong to the class of “Generalized Linear Models” (GLM)

Some Thoughts on Logistic and Softmax Classification

- Both logistic and softmax classification are discriminative model for linear classification
- These are special cases of the more general form of discriminative models for classification

$$p(y|\mathbf{x}) = \text{Bernoulli}(f(\mathbf{x})) \quad \text{OR} \quad p(y|\mathbf{x}) = \text{multinoulli}(f(\mathbf{x}))$$

- Logistic and softmax also belong to the class of “Generalized Linear Models” (GLM)
 - GLM can model outputs of various types, e.g., real, categorical, counts, positive reals, etc.

Some Thoughts on Logistic and Softmax Classification

- Both logistic and softmax classification are discriminative model for linear classification
- These are special cases of the more general form of discriminative models for classification

$$p(y|\mathbf{x}) = \text{Bernoulli}(f(\mathbf{x})) \quad \text{OR} \quad p(y|\mathbf{x}) = \text{multinoulli}(f(\mathbf{x}))$$

- Logistic and softmax also belong to the class of “Generalized Linear Models” (GLM)
 - GLM can model outputs of various types, e.g., real, categorical, counts, positive reals, etc.
 - Will look at GLM when we discuss exponential families

Some Thoughts on Logistic and Softmax Classification

- Both logistic and softmax classification are discriminative model for linear classification
- These are special cases of the more general form of discriminative models for classification

$$p(y|\mathbf{x}) = \text{Bernoulli}(f(\mathbf{x})) \quad \text{OR} \quad p(y|\mathbf{x}) = \text{multinoulli}(f(\mathbf{x}))$$

- Logistic and softmax also belong to the class of “Generalized Linear Models” (GLM)
 - GLM can model outputs of various types, e.g., real, categorical, counts, positive reals, etc.
 - Will look at GLM when we discuss exponential families
- Logistic/softmax can also be extended to handle nonlinear classification, e.g.,

Some Thoughts on Logistic and Softmax Classification

- Both logistic and softmax classification are discriminative model for linear classification
- These are special cases of the more general form of discriminative models for classification

$$p(y|\mathbf{x}) = \text{Bernoulli}(f(\mathbf{x})) \quad \text{OR} \quad p(y|\mathbf{x}) = \text{multinoulli}(f(\mathbf{x}))$$

- Logistic and softmax also belong to the class of “Generalized Linear Models” (GLM)
 - GLM can model outputs of various types, e.g., real, categorical, counts, positive reals, etc.
 - Will look at GLM when we discuss exponential families
- Logistic/softmax can also be extended to handle nonlinear classification, e.g.,
 - Using a deep neural network + sigmoid/softmax to model $f(x)$

Some Thoughts on Logistic and Softmax Classification

- Both logistic and softmax classification are discriminative model for linear classification
- These are special cases of the more general form of discriminative models for classification

$$p(y|\mathbf{x}) = \text{Bernoulli}(f(\mathbf{x})) \quad \text{OR} \quad p(y|\mathbf{x}) = \text{multinoulli}(f(\mathbf{x}))$$

- Logistic and softmax also belong to the class of “Generalized Linear Models” (GLM)
 - GLM can model outputs of various types, e.g., real, categorical, counts, positive reals, etc.
 - Will look at GLM when we discuss exponential families
- Logistic/softmax can also be extended to handle nonlinear classification, e.g.,
 - Using a deep neural network + sigmoid/softmax to model $f(x)$
 - Using Gaussian Processes (GP) which are essentially Bayesian kernel methods

Some Thoughts on Logistic and Softmax Classification

- Both logistic and softmax classification are discriminative model for linear classification
- These are special cases of the more general form of discriminative models for classification

$$p(y|\mathbf{x}) = \text{Bernoulli}(f(\mathbf{x})) \quad \text{OR} \quad p(y|\mathbf{x}) = \text{multinoulli}(f(\mathbf{x}))$$

- Logistic and softmax also belong to the class of “Generalized Linear Models” (GLM)
 - GLM can model outputs of various types, e.g., real, categorical, counts, positive reals, etc.
 - Will look at GLM when we discuss exponential families
- Logistic/softmax can also be extended to handle nonlinear classification, e.g.,
 - Using a deep neural network + sigmoid/softmax to model $f(x)$
 - Using Gaussian Processes (GP) which are essentially Bayesian kernel methods
 - We will look at such nonlinear classification models later..

Some Thoughts on Generative vs Discriminative Classification

- **Number of parameters:** Discriminative models have **fewer parameters** to be learned (e.g., just the weight vector/matrix \mathbf{w}/\mathbf{W} in case of logistic/softmax classification)

Some Thoughts on Generative vs Discriminative Classification

- **Number of parameters:** Discriminative models have **fewer parameters** to be learned (e.g., just the weight vector/matrix \mathbf{w}/\mathbf{W} in case of logistic/softmax classification)
- **Ease of parameter estimation:** Debatable as to which one is easier

Some Thoughts on Generative vs Discriminative Classification

- **Number of parameters:** Discriminative models have **fewer parameters** to be learned (e.g., just the weight vector/matrix \mathbf{w}/\mathbf{W} in case of logistic/softmax classification)
- **Ease of parameter estimation:** Debatable as to which one is easier
 - For “simple” class-conditionals (e.g., naïve Bayes assumption, diagonal covariance for Gaussian $p(\mathbf{x}|y)$), it is easier for generative classification model (often closed-form solution)

Some Thoughts on Generative vs Discriminative Classification

- **Number of parameters:** Discriminative models have **fewer parameters** to be learned (e.g., just the weight vector/matrix \mathbf{w}/\mathbf{W} in case of logistic/softmax classification)
- **Ease of parameter estimation:** Debatable as to which one is easier
 - For “simple” class-conditionals (e.g., naïve Bayes assumption, diagonal covariance for Gaussian $p(\mathbf{x}|y)$), it is easier for generative classification model (often closed-form solution)
 - Parameter estimation for discriminative models (logistic/softmax) usually requires iterative methods (although objective functions usually have global optima)

Some Thoughts on Generative vs Discriminative Classification

- **Number of parameters:** Discriminative models have **fewer parameters** to be learned (e.g., just the weight vector/matrix \mathbf{w}/\mathbf{W} in case of logistic/softmax classification)
- **Ease of parameter estimation:** Debatable as to which one is easier
 - For “simple” class-conditionals (e.g., naïve Bayes assumption, diagonal covariance for Gaussian $p(\mathbf{x}|\mathbf{y})$), it is easier for generative classification model (often closed-form solution)
 - Parameter estimation for discriminative models (logistic/softmax) usually requires iterative methods (although objective functions usually have global optima)
- **Dealing with missing features in the inputs:** Generative models can handle this easily (e.g., by integrating out the missing features while estimating the parameters)

Some Thoughts on Generative vs Discriminative Classification

- **Number of parameters:** Discriminative models have **fewer parameters** to be learned (e.g., just the weight vector/matrix \mathbf{w}/\mathbf{W} in case of logistic/softmax classification)
- **Ease of parameter estimation:** Debatable as to which one is easier
 - For “simple” class-conditionals (e.g., naïve Bayes assumption, diagonal covariance for Gaussian $p(\mathbf{x}|y)$), it is easier for generative classification model (often closed-form solution)
 - Parameter estimation for discriminative models (logistic/softmax) usually requires iterative methods (although objective functions usually have global optima)
- **Dealing with missing features in the inputs:** Generative models can handle this easily (e.g., by integrating out the missing features while estimating the parameters)
- **Inputs with features having mixed types:** Naturally handled by a generative model using an appropriate $p(\mathbf{x}|y)$ for each type of feature in the input. Difficult for discriminative models

Some Thoughts on Generative vs Discriminative Classification

- **Leveraging unlabeled data**: Generative models can handle this easily by treating the missing labels as latent variables and ideal for **Semi-supervised Learning**. Discriminative models can't do it easily

Some Thoughts on Generative vs Discriminative Classification

- **Leveraging unlabeled data:** Generative models can handle this easily by treating the missing labels as latent variables and ideal for **Semi-supervised Learning**. Discriminative models can't do it easily
- **Adding data from new classes:** Discriminative model will need to be re-trained. Generative model will just require estimating the class-conditionals of newly added classes

Some Thoughts on Generative vs Discriminative Classification

- **Leveraging unlabeled data:** Generative models can handle this easily by treating the missing labels as latent variables and ideal for **Semi-supervised Learning**. Discriminative models can't do it easily
- **Adding data from new classes:** Discriminative model will need to be re-trained. Generative model will just require estimating the class-conditionals of newly added classes
- **Have lots of labeled data:** Discriminative models usually work very well

Some Thoughts on Generative vs Discriminative Classification

- **Leveraging unlabeled data:** Generative models can handle this easily by treating the missing labels as latent variables and ideal for **Semi-supervised Learning**. Discriminative models can't do it easily
- **Adding data from new classes:** Discriminative model will need to be re-trained. Generative model will just require estimating the class-conditionals of newly added classes
- **Have lots of labeled data:** Discriminative models usually work very well
- **Final Verdict?** Despite generative classification having some clear advantages, both methods can be equally powerful (the actual choice may be dictated by the problem)

Some Thoughts on Generative vs Discriminative Classification

- **Leveraging unlabeled data:** Generative models can handle this easily by treating the missing labels as latent variables and ideal for **Semi-supervised Learning**. Discriminative models can't do it easily
- **Adding data from new classes:** Discriminative model will need to be re-trained. Generative model will just require estimating the class-conditionals of newly added classes
- **Have lots of labeled data:** Discriminative models usually work very well
- **Final Verdict?** Despite generative classification having some clear advantages, both methods can be equally powerful (the actual choice may be dictated by the problem)
 - Important to be aware of their strengths/weaknesses, and also the connections between these models

Some Thoughts on Generative vs Discriminative Classification

- **Leveraging unlabeled data:** Generative models can handle this easily by treating the missing labels as latent variables and ideal for **Semi-supervised Learning**. Discriminative models can't do it easily
- **Adding data from new classes:** Discriminative model will need to be re-trained. Generative model will just require estimating the class-conditionals of newly added classes
- **Have lots of labeled data:** Discriminative models usually work very well
- **Final Verdict?** Despite generative classification having some clear advantages, both methods can be equally powerful (the actual choice may be dictated by the problem)
 - Important to be aware of their strengths/weaknesses, and also the connections between these models
- **Possibility of a Hybrid Design?** Yes, Generative and Discriminative models also be combined, e.g.,

Some Thoughts on Generative vs Discriminative Classification

- **Leveraging unlabeled data**: Generative models can handle this easily by treating the missing labels as latent variables and ideal for **Semi-supervised Learning**. Discriminative models can't do it easily
- **Adding data from new classes**: Discriminative model will need to be re-trained. Generative model will just require estimating the class-conditionals of newly added classes
- **Have lots of labeled data**: Discriminative models usually work very well
- **Final Verdict?** Despite generative classification having some clear advantages, both methods can be equally powerful (the actual choice may be dictated by the problem)
 - Important to be aware of their strengths/weaknesses, and also the connections between these models
- **Possibility of a Hybrid Design**? Yes, Generative and Discriminative models also be combined, e.g.,
 - "Principled Hybrids of Generative and Discriminative Models" (Lassere et al, 2006)
 - "Deep Hybrid Models: Bridging Discriminative & Generative Approaches" (Kuleshov & Ermon, 2017)