

Linear Classification

Lili MOU

moull12@sei.pku.edu.cn

<http://sei.pku.edu.cn/~moull12>

23 April 2015

Outline

Introduction

Discriminant Functions

Probabilistic Generative Models

Probabilistic Discriminative Models

Bayesian Logistic Regression

Outline

Introduction

Discriminant Functions

Probabilistic Generative Models

Probabilistic Discriminative Models

Bayesian Logistic Regression

Prologue

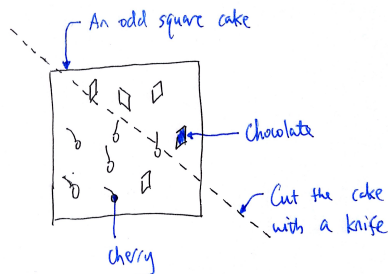
Linear classification is easy—my good old friend, logistic regression, always serves as a baseline method in various applications. Through a systematic study, however, we can grasp the main idea behind a range of machine learning techniques. This seminar also precedes our future discussion on GP classification.

References:

The materials basically follow Chapter 4 in *Pattern Recognition and Machine Learning*. Figures are taken (w/ or w/o modification) without further acknowledgment.

See also A. Webb, et al., *Statistical Pattern Recognition*.

The Linear Classification Problem



The Road Map

- ▶ Discriminant functions
- ▶ Probabilistic generative models
- ▶ Probabilistic discriminative models
- ▶ Bayesian logistic regression

On Non-Linearity

Generalized linearity: non-linear transformation by basis functions

- ▶ Feature engineering/selection

Kernels: transforming to a Hilbert space, fully characterized by a predefined “inner-product” operation

- ▶ SVM (a discriminant function)
- ▶ Gaussian processes (non-parametric Bayes)
- ▶ k -NN (slacked similarity measure, non-parametric discriminant)

Composition of non-linear activation functions

- ▶ Neural networks (finite feature learning, equivalent to learning a sophisticated kernel mapping input data to another Euclidean space)

My questions

- ▶ Can neural networks composite infinite dimensional features (with kernels)?
- ▶ Sparse kernel network?
- ▶ Or is it necessary? Or indeed, neural networks ARE non-parametric!

Outline

Introduction

Discriminant Functions

Probabilistic Generative Models

Probabilistic Discriminative Models

Bayesian Logistic Regression

The Decision Boundary

Let $\mathbf{x} = (x_1, \dots, x_n)^T$ be a data sample in \mathbb{R}^n

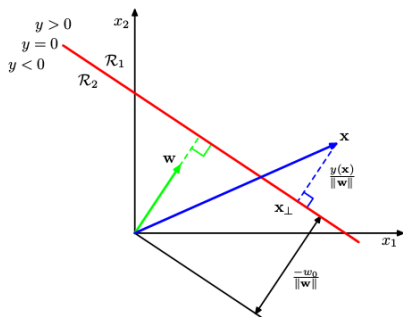
Let $y \in \{0, 1\}$ be the label of \mathbf{x} (a two-class problem)

A linear classifier always takes the form

$$y = \begin{cases} 1, & \text{if } \mathbf{w}^T \mathbf{x} + w_0 > 0 \\ 0, & \text{otherwise} \end{cases}$$

The *decision boundary* is a hyperplane in \mathbb{R}^{n-1}

$$\mathbf{w}^T \mathbf{x} + w_0 = 0$$

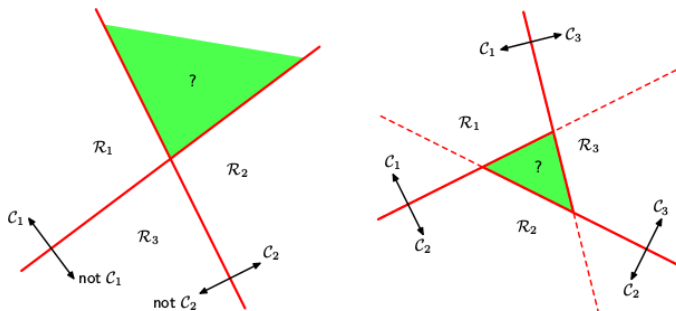


Multiclass problem

Let K be the number of classes

- ▶ *one-versus-the-rest* ($K - 1$ classifiers)
- ▶ *one-versus-one* + voting ($K(K - 1)/2$ classifiers)
- ▶ Scoring

$$y_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k0}$$



Linear Regression as Classification

Let the target value be 1-of- K coded. (K : # of classes)

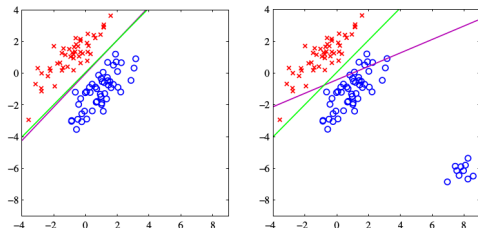
The scoring function is

$$y_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k0}$$

The least square objective function

$$J = \sum_{i=1}^m \sum_{k=1}^K \left(y_k(\mathbf{x}^{(i)}) - t^{(i)} \right)^2$$

The problem of being “too correct”—the target value is too far away from a Gaussian distribution.



Perception Learning Algorithm

Non-linear hard squashing function

$$y = \begin{cases} 1, & \text{if } \mathbf{w}^T \mathbf{x} + w_0 > 0 \\ -1, & \text{otherwise} \end{cases}$$

Perception learning algorithm

For each misclassified sample $x^{(i)}$:

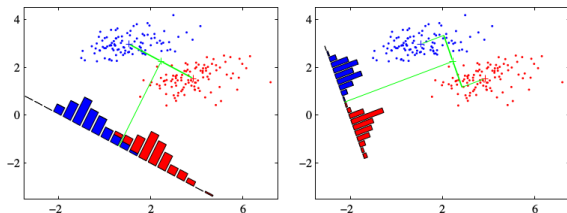
$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} + \eta \mathbf{x}^{(i)} t^{(i)}$$

For its convergence theorem, please refer to *Neural Networks and Machine Learning*.

Pros and Cons

- + Cares about mis-classified samples only
- Does not work with linearly inseparable data
- Does not generalize to multi-class problems
- + Introduces non-linear activation function

Fisher's Linear Discriminant



Heuristics:

- ▶ Attempt #1: Maximize the class separation
- ▶ Attempt #2 (the Fisher criterion):
Maximize the ratio of the between-class variance to the within-class variance (after projecting to a one-dimensional space, perpendicular to the decision boundary)

Formalization of Fisher's Linear Discriminant

Between-class variance

$$\mathbf{m}_1 = \frac{1}{N_1} \sum_{i \in \mathcal{C}_1} \mathbf{x}^{(i)}, \quad \mathbf{m}_2 = \frac{1}{N_2} \sum_{i \in \mathcal{C}_2} \mathbf{x}^{(i)}$$

Within-class variance

$$s_k^2 = \sum_{i \in \mathcal{C}_k} \left(y^{(i)} - m_k \right)^2$$

where

$$y_n = \mathbf{w}^T \mathbf{x}_n$$

The cost function

$$J(\mathbf{w}) = \frac{\mathbf{w}^T (\mathbf{m}_2 - \mathbf{m}_1)}{s_1^2 + s_2^2}$$

Solution to Fisher's Linear Discriminant

$$J(\mathbf{w}) = \frac{\mathbf{w}^T(\mathbf{m}_2 - \mathbf{m}_1)}{s_1^2 + s_2^2} = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}}$$

where

$$\mathbf{S}_B = (\mathbf{m}_2 - \mathbf{m}_1)(\mathbf{m}_2 - \mathbf{m}_1)^T$$

$$\mathbf{S}_W = \sum_{i \in \mathcal{C}_1} (\mathbf{x}_i - \mathbf{m}_1)(\mathbf{x}_i - \mathbf{m}_1)^T + \sum_{i \in \mathcal{C}_2} (\mathbf{x}_i - \mathbf{m}_2)(\mathbf{x}_i - \mathbf{m}_2)^T$$

Differentiating the cost function and setting it to 0, we obtain

$$(\mathbf{w} \mathbf{S}_B \mathbf{w}) \mathbf{S}_W \mathbf{w} = (\mathbf{w}^T \mathbf{S}_W \mathbf{w}) \mathbf{S}_B \mathbf{w}$$

Thus

$$\mathbf{w} \propto \mathbf{S}_W^{-1}(\mathbf{m}_2 - \mathbf{m}_1)$$

See *Pattern Recognition and Machine Learning* for multi-class scenarios.

Outline

Introduction

Discriminant Functions

Probabilistic Generative Models

Probabilistic Discriminative Models

Bayesian Logistic Regression

A Probabilistic Model

A data (x, y) is generated according to the following story

Step 1: Choose a box \mathcal{C}_k , i.e., $t = k$, with probability $p(\mathcal{C}_k)$

Step 2: Draw a ticket x from box k with probability $p(x|\mathcal{C}_k)$

The Bayesian network

$$(t) \longrightarrow (x)$$

Posterior probability

$$\begin{aligned} p(\mathcal{C}_1|\mathbf{x}) &= \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1) + p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)} \\ &= \frac{1}{1 + \exp(-a)} \triangleq \sigma(a) \end{aligned}$$

where

$$a = \ln \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)} = \ln \frac{p(\mathcal{C}_1|\mathbf{x})}{p(\mathcal{C}_2|\mathbf{x})}$$

Multi-class Scenarios

Softmax

$$p(\mathcal{C}_k|\mathbf{x}) = \frac{p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k)}{\sum_j p(\mathbf{x}|\mathcal{C}_j)p(\mathcal{C}_j)}$$

Predicting the class label

$$\hat{t} = \operatorname{argmax}_k p(\mathcal{C}_k|\mathbf{x})$$

Generative Models v.s. Discriminative Models

Training objectives

$$\underset{\Theta}{\text{maximize}} p(\mathbf{x}, \mathbf{t}; \Theta)$$

i.e.,

$$\underset{\Theta}{\text{maximize}} p(\mathbf{x}|\mathbf{t}; \Theta)p(\mathbf{t}; \Theta)$$

v.s.

$$\underset{\Theta}{\text{maximize}} p(\mathbf{t}|\mathbf{x}; \Theta)$$

Gaussian Inputs

Assumption

$$p(\mathbf{x}|\mathcal{C}_k) = \frac{1}{\sqrt{(2\pi)^n |\boldsymbol{\Sigma}|}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) \right\}$$

Consider a binary classification problem

$$p(\mathcal{C}_1|\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0)$$

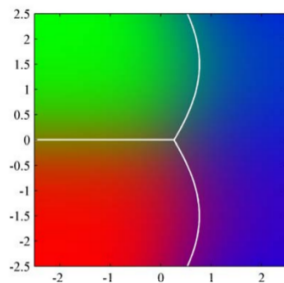
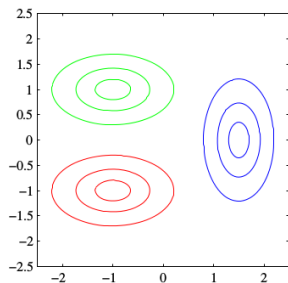
where

$$\begin{aligned} \mathbf{w} &= \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \\ w_0 &= -\frac{1}{2} \boldsymbol{\mu}_1^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1 + \frac{1}{2} \boldsymbol{\mu}_2^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_2 + \ln \frac{p(\mathcal{C}_1)}{p(\mathcal{C}_2)} \end{aligned}$$

Parameter estimation: Maximum likelihood estimation

See also Ch 2.3, 2.4 in *Statistical Pattern Recognition*.

Gaussian Classifiers



Shared covariance matrix $\Sigma \Rightarrow$ Linear decision boundary

Different $\Sigma_k \Rightarrow$ Quadratic decision boundary

Warning: Don't use Gaussian Classifiers

Discrete Features

Consider binary features $\mathbf{x} = (x_1, \dots, x_n)^T$, where $x_i \in \{0, 1\}$.

$p(\mathbf{x}|\mathcal{C}_k)$ has $2^n - 1$ independent free parameters

Naïve Bayes assumption: x_i independent

$$p(\mathbf{x}|\mathcal{C}_k) = \prod_{i=1}^n p(x_i|\mathcal{C}_k) = \prod_{i=1}^n \mu_{ki}^{x_i} (1 - \mu_{ki})^{1-x_i}$$

Posterior class distributions

$$p(\mathcal{C}_k|\mathbf{x}) = \sigma(a_k(\mathbf{x}))$$

$$a_k(\mathbf{x}) = \sum_{i=1}^n \{x_i \ln \mu_{ki} - (1 - x_i) \ln(1 - \mu_{ki})\} + \ln p(\mathcal{C}_k)$$

Linear in features \mathbf{x} !

Exponential Family

Let class-conditional distributions take the form

$$p(\mathbf{x}|\lambda_k) = h(\mathbf{x})g(\boldsymbol{\lambda}_k) \exp \{ \boldsymbol{\lambda}_k^T \mathbf{u}(\mathbf{x}) \}$$

Binary classification

$$a(\mathbf{x}) = (\boldsymbol{\lambda}_1 - \boldsymbol{\lambda}_2)^T \mathbf{x} + \ln g(\boldsymbol{\lambda}_1) - \ln g(\boldsymbol{\lambda}_2) + \ln p(\mathcal{C}_1) - \ln p(\mathcal{C}_2)$$

Multi-class problem

$$a_k(\mathbf{x}) = \boldsymbol{\lambda}_k^T \mathbf{x} + \ln g(\boldsymbol{\lambda}_k) + \ln p(\mathcal{C}_k)$$

Examples: Normal, exponential, log-normal, gamma, chi-squared, beta, Dirichlet, Bernoulli, categorical, Poisson, geometric, inverse Gaussian, von Mises, and von Mises-Fisher. Provided some parameters are fixed, Pareto, binomial, multinomial, and negative binomial are also in the exponential family.

Outline

Introduction

Discriminant Functions

Probabilistic Generative Models

Probabilistic Discriminative Models

Bayesian Logistic Regression

Logistic Regression

Under rather general assumptions, the posterior class distribution takes the form

$$y(\mathbf{x}) \triangleq p(\mathcal{C}_1|\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0)$$

The idea is to optimize directly the above posterior distribution.

⇒ Fewer parameters, better performance.

Consider a binary classification problem with n -dimensional feature space

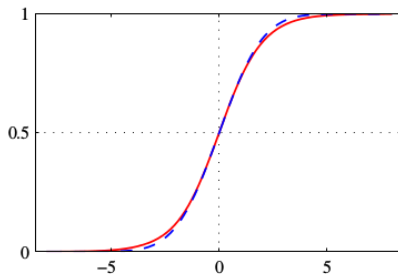
- ▶ Logistic regression: n parameters
- ▶ A Gaussian classifier:
 - + Means: $2n$
 - + Covariance matrices (shared): $n(n + 1)/2$
 - + Class prior: 1

Probit Regression

What if we substitute σ with other squashing functions?

Probit function

$$\Phi(a) = \int_{-\infty}^a \mathcal{N}(\theta|0, 1) d\theta$$



- ▶ Similar performance compared with logistic regression
- ▶ More prone to outliers (Why? Consider the decay rate \cup)
- ▶ Useful later

Outline

Introduction

Discriminant Functions

Probabilistic Generative Models

Probabilistic Discriminative Models

Bayesian Logistic Regression

Bayesian Learning

Let $\mathbf{y} \in \{0, 1\}^m$ denote the labels of training data ϕ_1, \dots, ϕ_m

Prior $p(\mathbf{w})$, which is a \$64,000,000 question

Likelihood

$$p(\mathbf{y}|\mathbf{w}) = \prod_{i=1}^m p(y^{(i)}|\mathbf{w}) = \prod_{i=1}^m \sigma(\mathbf{w}^T \phi^{(i)})^{t^{(i)}} (1 - \sigma(\mathbf{w}^T \phi^{(i)}))^{1-t^{(i)}}$$

Posterior

$$p(\mathbf{w}|\mathbf{y}) = \frac{p(\mathbf{w})p(\mathbf{y}|\mathbf{w})}{p(\mathbf{y})} \propto_{\mathbf{w}} p(\mathbf{w})p(\mathbf{y}|\mathbf{w}) = p(\mathbf{w}) \prod_{i=1}^m \sigma(\cdot)^{t^{(i)}} (1 - \sigma(\cdot))^{1-t^{(i)}}$$

Predictive density

$$p(y_*|\mathbf{y}) = \int d\mathbf{w} p(y_*|\mathbf{w}) \cdot p(\mathbf{w}|\mathbf{y})$$
$$\propto_{y_*} \int d\mathbf{w} \sigma(\mathbf{w}^T \phi_*) \cdot p(\mathbf{w}) \prod_{i=1}^m \sigma(\cdot)^{t^{(i)}} (1 - \sigma(\cdot))^{1-t^{(i)}}$$

Intractability

- ▶ Posterior is intractable due to the normalizing factor.
- ▶ Predictive density is intractable due to the integral.

We have to resort to approximations

- ▶ Sampling methods
 - Stochastic, usually asymptotically correct, hard to scale
- ▶ Deterministic methods
 - “Do things wrongly and hope they work”

Laplace Approximation

- ▶ Fit a Gaussian at a mode
- ▶ The standard deviation is chosen such that ...
the second-order derivative of the log probability matches
- ☺ The first-order derivative is always 0 at a mode
- ☺ Scale free in representing the unnormalized measure
- ☹ Real variables only
- ☹ Only local properties captured, multi-mode distributions?

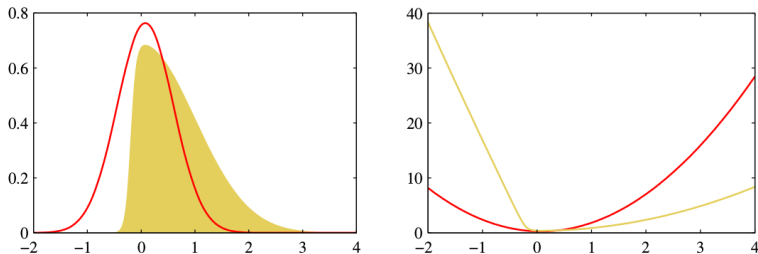


Figure 4.14 Illustration of the Laplace approximation applied to the distribution $p(z) \propto \exp(-z^2/2)\sigma(20z+4)$ where $\sigma(z)$ is the logistic sigmoid function defined by $\sigma(z) = (1 + e^{-z})^{-1}$. The left plot shows the normalized distribution $p(z)$ in yellow, together with the Laplace approximation centred on the mode z_0 of $p(z)$ in red. The right plot shows the negative logarithms of the corresponding curves.

Fitting a Gaussian

Let $p(z) = \frac{1}{Z} f(z)$ be a true distribution, where $Z = \int f(z) dz$

Step 1: Find a mode z_0 of $p(z)$, by gradient methods, say, satisfying

$$\left. \frac{df(z)}{dz} \right|_{z=z_0} = 0$$

Step 2: Consider a Taylor expansion of $\ln f(z)$ at z_0

$$\ln f(z) \simeq \ln f(z_0) - \frac{1}{2} A (z - z_0)^2$$

where A is given by

$$A = - \left. \frac{d^2}{dz^2} \ln f(z) \right|_{z=z_0}$$

Taking the exponential,

$$f(z) \simeq f(z_0) \exp \left\{ -\frac{A}{2} (z - z_0)^2 \right\}$$

Step 3: Normalize to a Gaussian distribution

$$q(z) = \left(\frac{A}{2\pi} \right)^{1/2} \exp \left\{ -\frac{A}{2} (z - z_0)^2 \right\}$$

Laplace Approximation for Multivariate Distributions

To approximate $p(\mathbf{z}) = \frac{1}{Z}f(\mathbf{z})$, where $\mathbf{z} \in \mathbb{R}^m$

We expand at mode \mathbf{z}_0

$$\ln f(\mathbf{z}) \simeq \ln f(\mathbf{z}_0) - \frac{1}{2}(\mathbf{z} - \mathbf{z}_0)^T \mathbf{A}(\mathbf{z} - \mathbf{z}_0)$$

where $\mathbf{A} = -\nabla\nabla \ln f(\mathbf{z})\Big|_{\mathbf{z}=\mathbf{z}_0}$, Hessian of $\ln f(\mathbf{z})$, serving as the *precision matrix* in a Gaussian distribution

Taking the exponential, we obtain

$$f(\mathbf{z}) \simeq f(\mathbf{z}_0) \exp \left\{ -\frac{1}{2}(\mathbf{z} - \mathbf{z}_0)^T \mathbf{A}(\mathbf{z} - \mathbf{z}_0) \right\}$$

Normalize it as a distribution, and then we have

$$q(\mathbf{z}) = \frac{|\mathbf{A}|^{1/2}}{(2\pi)^{m/2}} \exp \left\{ -\frac{1}{2}(\mathbf{z} - \mathbf{z}_0)^T \mathbf{A}(\mathbf{z} - \mathbf{z}_0) \right\} = \mathcal{N}(\mathbf{z}|\mathbf{z}_0, \mathbf{A}^{-1})$$

Bayesian Logistic Regression: A Revisit in Earnest

Prior: Gaussian, which is natural¹

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_0, \mathbf{S}_0)$$

Posterior: $p(\mathbf{w} | \mathbf{t}) \propto p(\mathbf{w})p(\mathbf{t} | \mathbf{w})$

$$\ln p(\mathbf{w} | \mathbf{t}) = -\frac{1}{2}(\mathbf{w} - \mathbf{m}_0)^T \mathbf{S}_0^{-1}(\mathbf{w} - \mathbf{m}_0) \quad [\text{prior}]$$

$$+ \sum_{i=1}^m \left\{ t^{(i)} \ln y^{(i)} + (1 - t^{(i)}) \ln(1 - y^{(i)}) \right\} \quad [\text{likelihood}]$$

+ const

$$\mathbf{S}_N = -\nabla \nabla \ln p(\mathbf{w} | \mathbf{t}) = \mathbf{S}_0^{-1} + \sum_{i=1}^m y^{(i)}(1 - y^{(i)}) \phi_n \phi_n^T$$

Hence, the Laplace approximation to the posterior is

$$q(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{w}_{\text{MAP}}, \mathbf{S}_N)$$

¹Mathematicians always choose priors for the sake of convenience rather than approaching God.

Predictive Density

$$p(\mathcal{C}_1|\phi_*, \mathbf{t}) = \int p(\mathcal{C}_1|\phi_*, \mathbf{w})p(\mathbf{w}|\mathbf{t}) d\mathbf{w} \simeq \int \sigma(\mathbf{w}^T \phi_*)q(\mathbf{w}) d\mathbf{w}$$

$$p(\mathcal{C}_2|\phi_*, \mathbf{t}) = 1 - p(\mathcal{C}_1|\phi_*, \mathbf{t})$$

Plan

- ▶ Change it to a univariate integral
- ▶ Substitute sigmoid with a probit function, which is then convolved with a normal

$$\int \sigma(\cdot)\mathcal{N}(\cdot) d\cdot \simeq \int \Phi(\cdot)\mathcal{N}(\cdot) d\cdot = \Phi(\cdot) \simeq \sigma(\cdot)$$

The Dirac Delta Function

Let δ be the *Dirac delta* function, loosely thought of a function such that

- ▶ Gaussian distribution peaked at 0 with standard deviation $\rightarrow 0$

- ▶
$$\delta(x) = \begin{cases} +\infty, & \text{if } x = 0 \\ 0, & \text{otherwise} \end{cases}$$

δ function satisfies

$$\int \delta(a - x) f(a) da = f(x)$$

and specifically

$$\int \delta(a - x) da = 1$$

Deriving the Predictive Density

$$p(\mathcal{C}_1|\phi_*) \simeq \int \sigma(\mathbf{w}^T \phi_*) q(\mathbf{w}) d\mathbf{w} \quad \text{[Laplace approx.]}$$

$$\sigma(\mathbf{w}^T \phi) = \int \delta(a - \mathbf{w}^T \phi) \sigma(a) da \quad \text{[Def. of } \delta \text{]}$$

$$\begin{aligned} p(\mathcal{C}_1|\phi_*) &= \int \int \delta(a - \mathbf{w}^T \phi) \sigma(a) da q(\mathbf{w}) d\mathbf{w} \\ &= \int \int \sigma(a) \delta(a - \mathbf{w}^T \phi_*) q(\mathbf{w}) d\mathbf{w} da \\ &\stackrel{\Delta}{=} \int \sigma(a) p(a) da \end{aligned}$$

where

$$p(a) \stackrel{\Delta}{=} \int \delta(a - \mathbf{w}^T \phi_*) q(\mathbf{w}) d\mathbf{w}$$

We now argue that $p(a)$ is Gaussian

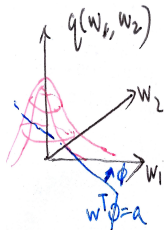
Deriving the Predictive Density (2)

$$\begin{aligned} p(a) &= \int \int \cdots \int \delta(a - \mathbf{w}^T \boldsymbol{\phi}_*) q(\mathbf{w}) d w_1 d w_2 \cdots d w_n \\ &= \int \int \cdots \int q(\tilde{\mathbf{w}}) d w_2 \cdots d w_n \end{aligned}$$

$\tilde{\mathbf{w}}$ is such that $\tilde{\mathbf{w}}^T \boldsymbol{\phi}_* = a$

We can also verify that

$$\begin{aligned} \int p(a) da &= \int \int \delta(a - \mathbf{w}^T \boldsymbol{\phi}_*) q(\mathbf{w}) d \mathbf{w} da \\ &= \int \int \delta(a - \mathbf{w}^T \boldsymbol{\phi}_*) q(\mathbf{w}) da d \mathbf{w} \\ &= \int q(\mathbf{w}) d \mathbf{w} \int \delta(a - \mathbf{w}^T \boldsymbol{\phi}_*) da \\ &= 1 \end{aligned}$$



Deriving the Predictive Density (3)

$$\mu_a = \mathbb{E}[a] = \int p(a)a \, da = \int q(\mathbf{w})\mathbf{w}^T \boldsymbol{\phi}_* \, d\mathbf{w} = \mathbf{w}_{\text{MAP}}^T \boldsymbol{\phi}_*$$

$$\begin{aligned}\sigma_a^2 &= \text{var}[a] = \int p(a) \{a^2 - \mathbb{E}[a]^2\} \, da \\ &= \int q(\mathbf{w}) \{(\mathbf{w}^T \boldsymbol{\phi})^2 - (\mathbf{m}_N^T \boldsymbol{\phi})^2\} \, d\mathbf{w} \\ &= \boldsymbol{\phi}^T \mathbf{S}_N \boldsymbol{\phi}\end{aligned}$$

Thus

$$\begin{aligned}p(\mathcal{C}_1 | \mathbf{t}) &\simeq \int \sigma(a)p(a) \, da = \int \sigma(a)\mathcal{N}(a|\mu_a, \sigma_a^2) \, da \\ &\simeq \int \Phi(\lambda a)\mathcal{N}(a|\mu_a, \sigma_a^2) \, da = \Phi\left(\frac{\mu_a}{(\lambda^{-2} + \sigma_a^2)^{1/2}}\right) \simeq \sigma(\kappa(\sigma_a^2)\mu_a)\end{aligned}$$

$\lambda = \sqrt{\pi/8}$, $\kappa(\sigma_a^2) = (1 + \pi\sigma_a^2/8)^{-1/2}$, chosen such that the rescaled probit function has the same slope as sigmoid at the origin.

$$\Phi(\cdot) * \phi(\cdot) = \Phi(\cdot)$$

Let $X \sim \mathcal{N}(a, b^2)$ and $Y \sim \mathcal{N}(c, d^2)$

$$\begin{aligned}\Pr\{X \leq Y\} &= \int_{-\infty}^{\infty} \Pr\{X \leq Y | Y = w\} \phi\left(\frac{w - c}{d}\right) dw \\ &= \int_{-\infty}^{\infty} \Pr\{X \leq w\} \phi\left(\frac{w - c}{d}\right) dw \\ &= \int_{-\infty}^{\infty} \Phi\left(\frac{w - a}{b}\right) \phi\left(\frac{w - c}{d}\right) dw\end{aligned}$$

By noticing that

$$X - Y \sim \mathcal{N}(a - c, b^2 + d^2)$$

We have

$$\begin{aligned}\Pr\{X \leq Y\} &= \Pr\{X - Y \leq 0\} \\ &= \Phi\left(\frac{-a + c}{\sqrt{b^2 + d^2}}\right)\end{aligned}$$

Take-Home Messages

- ▶ Discriminant functions
- ▶ Probabilistic generative models (don't use it)
- ▶ Probabilistic discriminative models
- ▶ Bayesian logistic regression
 - ✂ Sampling methods
 - 📦 Deterministic approximations
 - ▶ Laplace approximation (fitting a Gaussian at a mode)
 - ▶ Substitute the sigmoid function with a probit function
 - ▶ Analytical solutions
 - ▶ The decision boundary is the same with equal prior probabilities

Thanks for Listening!