# Logistic and Softmax Regression

CS771: Introduction to Machine Learning

Nisheeth

# Evaluation Measures for Regression Models

- Plotting the prediction $\hat{y}_n$ vs truth $y_n$ for the validation/test set
- Residual Sum of Squares (RSS) on the validation/test set

$$RSS(\boldsymbol{w}) = \sum_{n=1}^{N} (y_n - \hat{y}_n)^2$$

Plots of true vs predicted outputs and $R^2$ for two regression models



degree 1. R2 on Test = 0.473

- RMSE (Root Mean Squared Error) $\triangleq \sqrt{\dfrac{1}{N} RSS(w)}$

- Coefficient of determination or $R^2$

$$R^2 = 1 - \frac{\sum_{n=1}^{N}(y_n - \hat{y}_n)^2}{\sum_{n=1}^{N}(y_n - \bar{y})^2}$$

"relative" error w.r.t. a model that makes a constant prediction $\bar{y}$ for all inputs

Unlike RSS and RMSE, it is always between 0 and 1 and hence interpretable

$\bar{y}$ is empirical mean of true responses, i.e., $\frac{1}{N}\sum_{n=1}^{N} y_n$



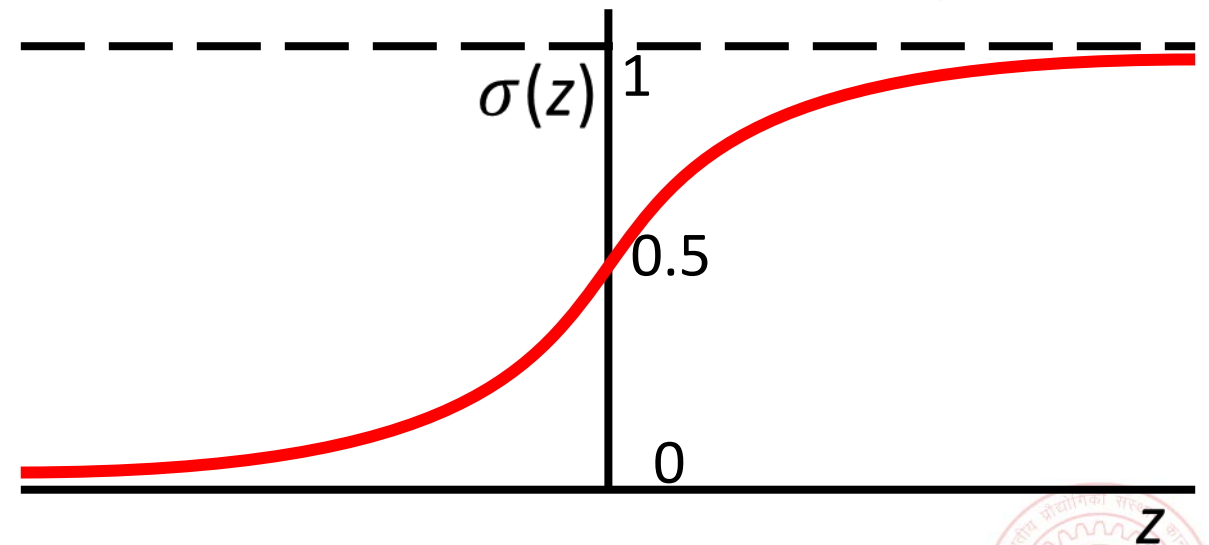degree 2. R2 on Test = 0.813

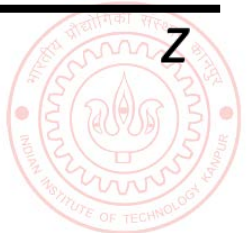Pic from MLAPP (Murphy)

CS771: Intro to ML

# Logistic Regression (LR)

- A probabilistic model for binary classification
- Learns the PMF of the output label given the input, i.e., $p(y|x)$
- A discriminative model: Does not model inputs $x$ (only relationship b/w $x$ and $y$)
- Uses the sigmoid function to define the conditional probability of $y$ being 1

$$\mu_x = p(y = 1 | w, x) = \sigma(w^\top x)$$

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

$$= \frac{\exp(w^\top x)}{1 + \exp(w^\top x)}$$

A linear model



- Here $w^\top x$ is the <u>score</u> for input $x$. The sigmoid turns it into a probability

# LR: Decision Boundary

- At the decision boundary where both classes are equiprobable

$$p(y = 1|x, w) \quad = \quad p(y = 0|x, w)$$

$$\frac{\exp(w^\top x)}{1 + \exp(w^\top x)} \quad = \quad \frac{1}{1 + \exp(w^\top x)}$$

$$\exp(w^\top x) \quad = \quad 1$$

$$w^\top x \quad = \quad 0$$

A linear hyperplane



$w^\top x = 0$

- Very large positive $w^\top x$ means $p(y = 1|w, x)$ close to 1
- Very large negative $w^\top x$ means $p(y = 0|w, x)$ close to 1
- At decision boundary, $w^\top x = 0$ implies $p(y = 1|w, x) = p(y = 0|w, x) = 0.5$

# MLE for Logistic Regression

- Likelihood (PMF of each input's label) is Bernoulli with prob $\mu_n = \dfrac{\exp(\boldsymbol{w}^\top \boldsymbol{x}_n)}{1+\exp(\boldsymbol{w}^\top \boldsymbol{x}_n)}$

Assumed 0/1, not -1/+1

$$p(y_n|\boldsymbol{w}, \boldsymbol{x}_n) = \text{Bernoulli}(\mu_n) = \mu_n^{y_n}(1-\mu_n)^{1-y_n}$$

- Overall likelihood, assuming i.i.d. observations

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

- The negative log-likelihood $NLL(\boldsymbol{w}) = -\log p(\boldsymbol{y}|\boldsymbol{w}, \boldsymbol{X})$ simplifies to

"cross-entropy" loss (a popular loss function for classification)

Loss function

$$NLL(\boldsymbol{w}) = \sum_{n=1}^{N} -[y_n \log \mu_n + (1-y_n)\log(1-\mu_n)]$$

Very large loss if $y_n$ close to 1 and $\mu_n$ close to 0, or vice-versa

- Plugging in $\mu_n = \dfrac{\exp(\boldsymbol{w}^\top \boldsymbol{x}_n)}{1+\exp(\boldsymbol{w}^\top \boldsymbol{x}_n)}$ and simplifying

Good news: For LR, NLL is convex

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

# An Alternate Notation

- If we assume the label $y_n$ as -1/+1 (not 0/1), the likelihood can be written as

$$p(y_n|\boldsymbol{w}, \boldsymbol{x}_n) = \frac{1}{1 + \exp(-y_n \boldsymbol{w}^\top \boldsymbol{x}_n)} = \sigma(y_n \boldsymbol{w}^\top \boldsymbol{x}_n)$$

- Slightly more convenient notation: A single expression gives the probabilities of both possible label values

- In this case, the total negative log-likelihood will be

$$NLL(\boldsymbol{w}) = \sum_{n=1}^{N} -\log p(y_n|\boldsymbol{w}, \boldsymbol{x}_n) = \sum_{n=1}^{N} \log\left(1 + \exp(-y_n \boldsymbol{w}^\top \boldsymbol{x}_n)\right)$$

# MAP Estimation for Logistic Regression

- Need a prior on the weight vector $\boldsymbol{w} \in \mathbb{R}^D$

- Just like probabilistic linear regression, can use a zero-mean Gaussian prior

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

Or NLL – log of prior

- The MAP objective (log of posterior) will be log-likelihood + log of prior

- Therefore the MAP solution (ignoring terms that don't depend on $\boldsymbol{w}$) will be

$$\widehat{\boldsymbol{w}}_{MAP} = \arg\min_{\boldsymbol{w}} NLL(\boldsymbol{w}) + \frac{\lambda}{2}\boldsymbol{w}^\top\boldsymbol{w}$$

Good news: convex objective

- Just like MLE case, no closed form solution. Iterative opt methods needed
  - Highly efficient solvers (both first and second order) exist for MLE/MAP estimation for LR
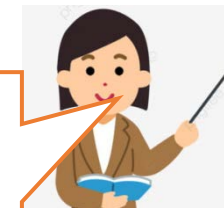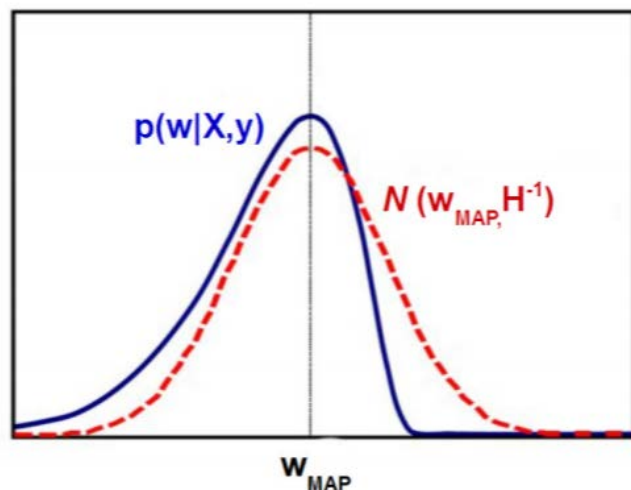
# Fully Bayesian Inference for Logistic Regression

- Doing fully Bayesian inference would require computing the posterior

Gaussian    Bernoulli

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

Unfortunately, Gaussian and Bernoulli are not conjugate with each other, so analytic expression for the posterior can't be obtained unlike prob. linear regression
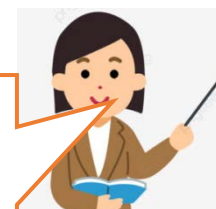
- Need to approximate the posterior in this case

- We will use a simple approximation called Laplace approximation



Approximates the posterior of **w** by a Gaussian whose mean is the MAP solution $\widehat{\boldsymbol{w}}_{MAP}$ and covariance matrix is the inverse of the Hessian (Hessian: second derivative of the negative log-posterior of the LR model)

Can also employ more advanced posterior approximation methods, like MCMC and variational inference
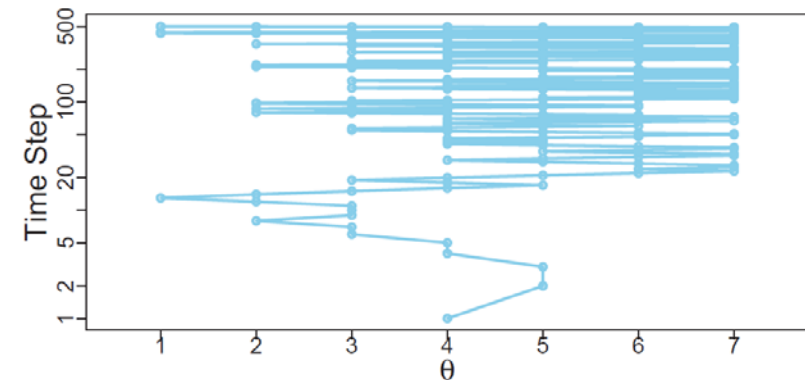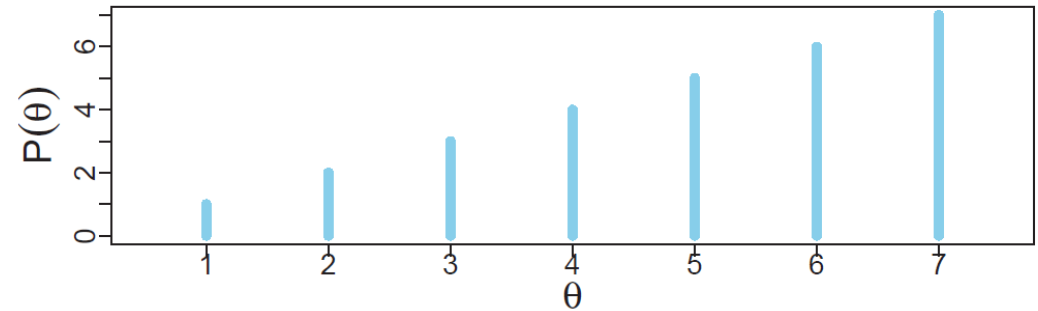
# Quick detour: MCMC sampling

- So far we have seen Bayesian calculations with conjugate priors, where the posterior is guaranteed to have the same form as the prior

- In many interesting situations in probabilistic ML, this will not be true

- MCMC-based approximation of the posterior is a godsend in such situations

- Let us build some intuition about MCMC sampling

- Suppose there is a greedy cat, who wanders from house to house expecting to be fed

# Cat food sampling: The Metropolis algorithm

- Assume that the cat only frequents a row of seven houses on a block

- Let the probability of the cat being fed at any of the houses be represented by the graph on the right

- The cat
  - Flips a coin to decide whether to go left or right
  - If she thinks the probability of being fed there is higher, she goes there
  - If not, then she goes there probabilistically as a function of the ratio of the probability of being fed in either of the two places

- Sampled long enough, the distribution of the cat's visits reflects the true underlying distribution of her being fed in the block

# Metropolis-Hastings for Bayesian inference

- We are interested in estimating the distribution of model parameters after seeing data p(θ|y)

- We sample possible proposal θ from a 'jump' distribution that conditions on the immediately previous θ value (say a Gaussian)

- Find the ratio of the proposed posterior distribution to the current one

$$r = \frac{f_{\underset{\sim}{\theta}}(\theta_i^*) f_{Y|\underset{\sim}{\theta}}(y|\theta_i^*)}{f_{\underset{\sim}{\theta}}(\theta_{i-1}) f_{Y|\underset{\sim}{\theta}}(y|\theta_{i-1})}$$
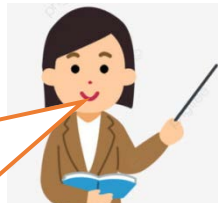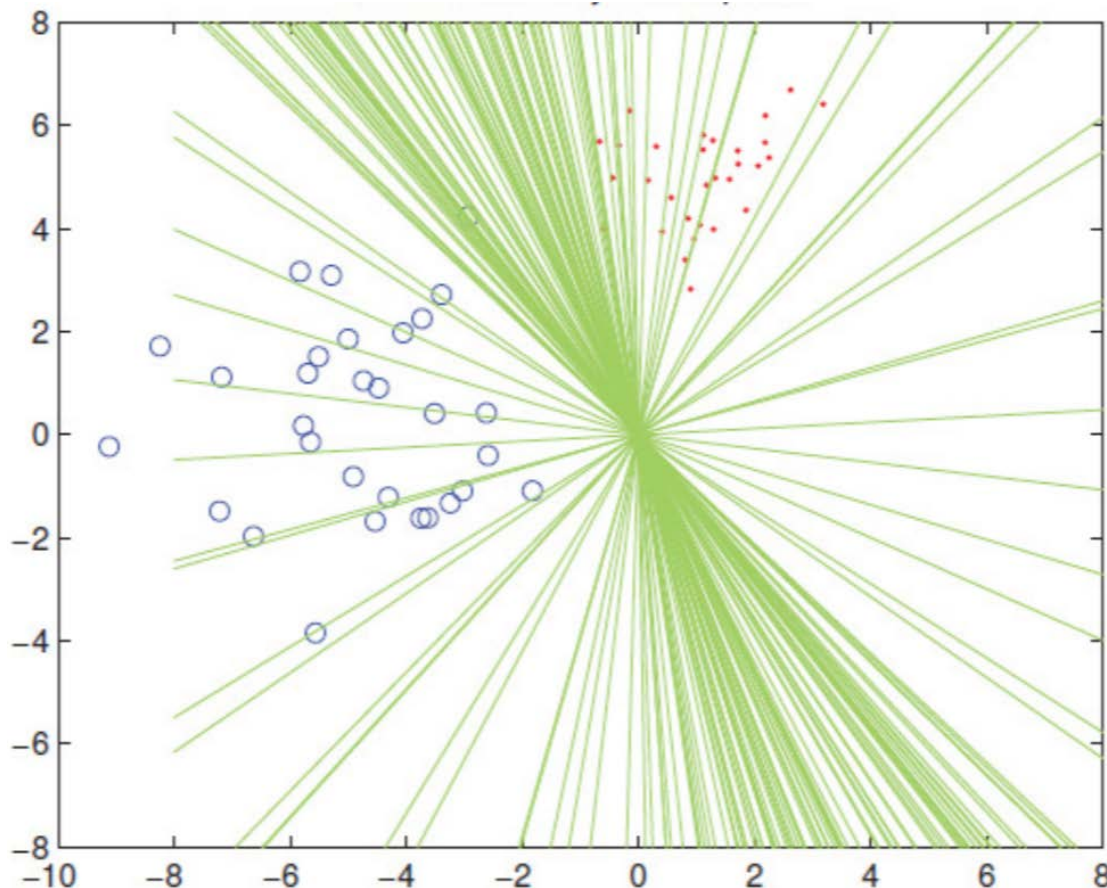
- Generate a random number a between 0 and 1

- If r > a, accept the proposed θ, otherwise stick with the earlier θ

- After tens of thousands of samples, the histogram of θ approximates p(θ|y)

- Incredibly simple but powerful way of approximating Bayesian posteriors

A nice tutorial for applying MCMC approximation for logistic regression

# Posterior for LR: An Illustration

- Can sample from the posterior of the LR model

- Each sample will give a weight vec defining a hyperplane separator



Not all separators are equally good; their goodness depends on their posterior probabilities

When making predictions, we can still use all of them but weighted by their importance based on their posterior probabilities

That's exactly what we do when computing the predictive distribution

# Logistic Regression: Predictive Distribution

- When using MLE/MAP solution $\widehat{\boldsymbol{w}}_{opt}$, can use the plug-in predictive distribution

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

$$\approx p(y_* = 1|\widehat{\boldsymbol{w}}_{opt}, \boldsymbol{x}_*) = \sigma(\widehat{\boldsymbol{w}}_{opt}^{\top}\boldsymbol{x}_n)$$

$$p(y_*|\boldsymbol{x}_*, \boldsymbol{X}, \boldsymbol{y}) = \text{Bernoulli}[\sigma(\widehat{\boldsymbol{w}}_{opt}^{\top}\boldsymbol{x}_n)]$$

- When using fully Bayesian inference, we must compute the posterior predictive

$$p(y_* = 1|\boldsymbol{x}_*, \boldsymbol{X}, \boldsymbol{y}) = \int p(y_* = 1|\boldsymbol{w}, \boldsymbol{x}_*)p(\boldsymbol{w}|\boldsymbol{X}, \boldsymbol{y})d\boldsymbol{w}$$
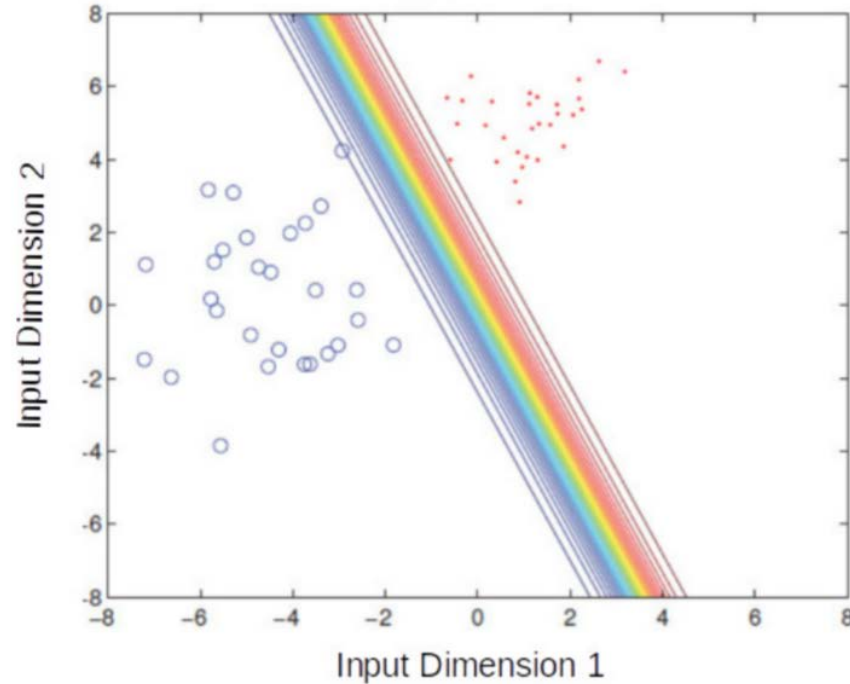
Integral not tractable and must be approximated

sigmoid

Gaussian (if using Laplace approx.)

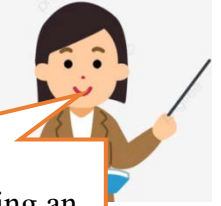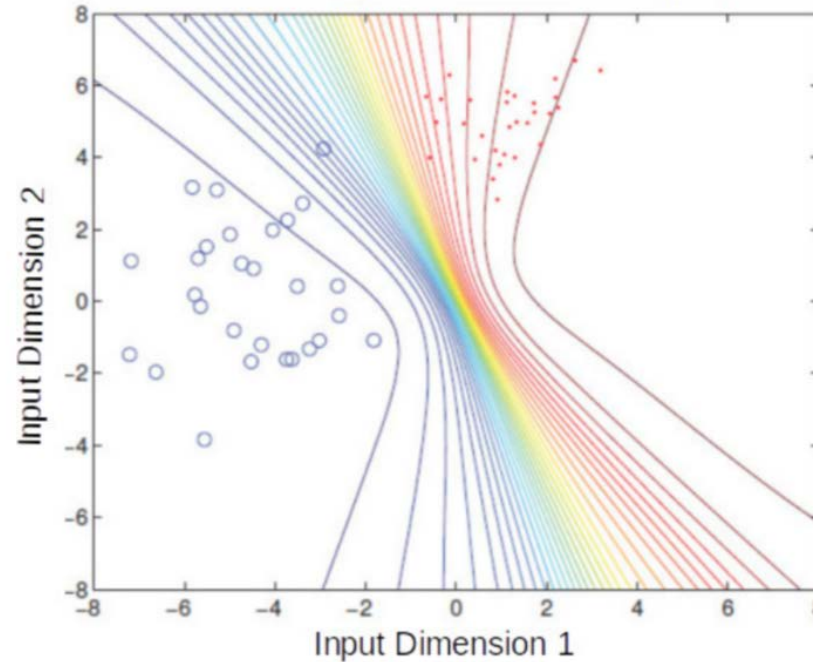Monte-Carlo approximation of this integral is one possible way

Generate $M$ samples $\boldsymbol{w}_1, \boldsymbol{w}_2, \ldots, \boldsymbol{w}_M$, from the Gaussian approx. of posterior and use $p(y_* = 1|\boldsymbol{x}_*, \boldsymbol{X}, \boldsymbol{y}) \approx \frac{1}{M}\sum_{m=1}^{M} p(y_* = 1|\boldsymbol{w}_m, \boldsymbol{x}_*) = \frac{1}{M}\sum_{m=1}^{M}\sigma(\boldsymbol{w}_m^{\top}\boldsymbol{x}_n)$

# LR: Plug-in Prediction vs Postrerior Averaging



Posterior averaging is like using an ensemble of models. In this example, each model is a linear classifier but the ensemble-like effect resulted in nonlinear boundaries

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

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

Softmax function

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

Also note that $\sum_{\ell=1}^{K} \mu_{n\ell} = 1$ for any input $\boldsymbol{x}_n$

- $K$ weight vecs $\boldsymbol{w}_1, \boldsymbol{w}_2, \ldots, \boldsymbol{w}_K$ (one per class), each $D$-dim, and $\boldsymbol{W} = [\boldsymbol{w}_1, \boldsymbol{w}_2, \ldots, \boldsymbol{w}_K]$
- Each likelihood $p(y_n | \boldsymbol{x}_n, \boldsymbol{W})$ is a multinoulli distribution. Therefore total likelihood

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

Notation: $y_{n\ell} = 1$ if true class of $\boldsymbol{x}_n$ is $\ell$ and $y_{n\ell'} = 0 \ \forall \ \ell' \neq \ell$

- Can do MLE/MAP/fully Bayesian estimation for $\boldsymbol{W}$ similar to LR model