

CPSC 340: Machine Learning and Data Mining

Multi-Class Regression

Fall 2016

Admin

- **Midterm:**
 - Grades/solutions will be posted later this week.
- **Assignment 4:**
 - Posted, due November 14.
- **Extra office hours:**
 - Thursdays from 4:30-5:30 in ICICS X836.

Last Time: L1-Regularization

- We discussed **L1-regularization**:

$$f(w) = \frac{1}{2} \|Xw - y\|^2 + \lambda \|w\|_1$$

- Also known as “LASSO” and “basis pursuit denoising”.
 - **Regularizes ‘w’** so we decrease our test error (like L2-regularization).
 - Yields **sparse ‘w’** so it selects features (like L0-regularization).
- **Properties:**
 - It’s **convex and fast** to minimize (proximal-gradient).
 - Solution is **not unique**.
 - Tends to yield **false positives**.

Extensions of L1-Regularization

- “Elastic net” uses L2-regularization plus L1-regularization.
 - Solution is still sparse but is now unique.
 - Slightly better with feature dependence: selects both “mom” and “mom2”.
- “Bolasso” runs L1-regularization on bootstrap samples.
 - Selects features that are non-zero in all samples.
 - Much less sensitive to false positives.
- There are *many* non-convex regularizers (square-root, “SCAD”):
 - Much less sensitive to false positives.
 - But computing global minimum is hard.

Last: Maximum Likelihood Estimation

- We discussed computing 'w' by **maximum likelihood estimation (MLE)**:

$$p(y | X, w)$$

- This is equivalent to minimizing **negative log-likelihood (NLL)**:

$$f(w) = - \sum_{i=1}^n \log(p(y_i | x_i, w))$$

- For **logistic regression**, probability is $w^T x_i$ passed through **sigmoid**.

$$p(y_i | x_i, w) = \frac{1}{1 + \exp(-y_i w^T x_i)}$$

- And MLE is minimum of:

$$f(w) = \sum_{i=1}^n \log(1 + \exp(-y_i w^T x_i))$$

- With regularization, similar to SVMs. But **gives probabilities**.

Maximum Likelihood and Least Squares

- Many of our objective functions can be written as an MLE.
- For example, consider **Gaussian likelihood** with mean of $w^T x_i$:

$$p(y_i | x_i, w) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(w^T x_i - y_i)^2}{2\sigma^2}\right)$$

- So the NLL is given by:

$$f(w) = -\sum_{i=1}^n \log(p(y_i | x_i, w))$$

$$= -\sum_{i=1}^n \log\left(\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(w^T x_i - y_i)^2}{2\sigma^2}\right)\right)$$

$$= -\sum_{i=1}^n \left[\underbrace{\log\left(\frac{1}{\sqrt{2\pi\sigma^2}}\right)}_{\text{constant in terms of 'w'}} + \underbrace{\log\left(\exp\left(-\frac{(w^T x_i - y_i)^2}{2\sigma^2}\right)\right)}_{\text{operations cancel}} \right]$$

$$= (\text{constant}) - \sum_{i=1}^n \left(-\frac{(w^T x_i - y_i)^2}{2\sigma^2}\right)$$

$$= (\text{constant}) + \frac{1}{2\sigma^2} \sum_{i=1}^n (w^T x_i - y_i)^2$$

positive constant
so doesn't change solution.

Maximum Likelihood and Least Squares

- Many of our objective functions can be written as an MLE.
- For example, consider **Gaussian likelihood** we mean of $w^T x_i$:

$$p(y_i | x_i, w) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(w^T x_i - y_i)^2}{2\sigma^2}\right)$$

- So we can minimize NLL by minimizing:

$$\begin{aligned} f(w) &= \frac{1}{2} \sum_{i=1}^n (w^T x_i - y_i)^2 \\ &= \frac{1}{2} \|Xw - y\|^2 \end{aligned}$$

- So **least squares is MLE under Gaussian likelihood**.
 - With a **Laplace likelihood** you would absolute error.

Problem with Maximum Likelihood Estimation

- Maximum likelihood estimate maximizes:

$$p(y | X, w)$$

- It's a bit weird:
 - “Find the ‘w’ that makes ‘y’ have the highest probability given ‘X’ and ‘w’.”
- A **problem with MLE**:
 - ‘y’ could be very likely for some **very unlikely ‘w’**.
 - E.g., complex model that overfit by memorizing the data.
- What we really want:
 - “Find the **‘w’ that has the highest probability** given ‘X’ and ‘y’.”

Maximum a Posteriori (MAP) Estimation

- Maximum a posteriori (MAP) maximizes what we want:

$$p(w | X, y)$$

"posterior" probability of 'w' given data $\{X, y\}$

- Using Bayes' rule, we have

$$p(w | X, y) = \frac{p(y | X, w) p(w | X)}{p(y | X)} \propto p(y | X, w) p(w | X)$$

$$= p(y | X, w) p(w)$$

"likelihood" "prior"

↓ Assume 'w' does not depend on 'X'.

Maximum a Posteriori (MAP) Estimation

- Maximum a posteriori (MAP) maximizes what we want:

$$p(w | X, y) \propto p(y | X, w) p(w)$$

"posterior" "likelihood" "prior"

- Prior $p(w)$ is 'belief' that 'w' is the correct before seeing data:
 - Can take into account that complex models can overfit.

- If we again minimize the negative of the logarithm, we get:

$$\begin{aligned} -\log(p(w | X, y)) &= -\log(p(y | X, w)) - \log(p(w)) + (\text{constant}) \\ &= -\sum_{i=1}^n \log(p(y_i | x_i, w)) - \sum_{j=1}^d \log(p(w_j)) + (\text{constant}) \end{aligned}$$

↓ IID data ↓ prior over the w_j are independent

MAP Estimation and Regularization

- While many losses are equivalent to NLLs, many **regularizers are equivalent to negative log-priors**.
- Assume each w_j comes from a Gaussian (0-mean, $1/\lambda$ variance):

$$p(w_j) = \frac{1}{\sqrt{2\pi(1/\lambda)}} \exp\left(-\frac{(w_j - 0)^2}{2(1/\lambda)}\right)$$

- Then the **log-prior** is:

$$\begin{aligned} \log(p(w_j)) &= -\log(\sqrt{2\pi(1/\lambda)}) + \log(\exp(-\frac{\lambda}{2} w_j^2)) \\ &= (\text{constant}) - \frac{\lambda}{2} w_j^2 \end{aligned}$$

- And **negative log-prior** over all 'j' is:

$$-\log(p(w)) = -\sum_{j=1}^d \log(p(w_j)) = (\text{constant}) + \sum_{j=1}^d \frac{\lambda}{2} w_j^2 = \frac{\lambda}{2} \|w\|^2 + \text{const.} \quad (L_2\text{-regularization})$$

MAP Estimation and Regularization

- MAP estimation gives **link between probabilities and loss functions**.
 - Gaussian likelihood and Gaussian prior gives L2-regularized least squares.

$$\text{If } p(y_i | x_i, w) \propto \exp\left(-\frac{(w^T x_i - y_i)^2}{2\sigma^2}\right) \quad p(w_j) \propto \exp\left(-\frac{\lambda}{2} w_j^2\right)$$

Then MAP estimation is equivalent to minimizing $f(w) = \frac{1}{2\sigma^2} \|Xw - y\|^2 + \frac{\lambda}{2} \|w\|^2$

- Sigmoid likelihood and Gaussian prior gives L2-regularized logistic regression:

$$\text{If } p(y_i | x_i, w) = \frac{1}{1 + \exp(-y_i w^T x_i)} \quad \text{and } p(w_j) \propto \exp\left(-\frac{\lambda}{2} w_j^2\right)$$

then MAP estimate is minimum of $f(w) = \sum_{i=1}^n \log(1 + \exp(-y_i w^T x_i)) + \frac{\lambda}{2} \|w\|^2$

As $n \rightarrow \infty$ effect of prior/regularizer goes to 0

↑ For MAP estimation the constant changes solution.

But setting $\sigma^2 \neq 1$ is equivalent to changing λ so we usually use $\sigma^2 = 1$

Why do we care about MLE and MAP?

- Unified way of thinking about many of our tricks?
 - Laplace smoothing in naïve Bayes can be viewed as regularization.
- Remember our two ways to **reduce complexity** of a model:
 - **Model averaging** (ensemble methods).
 - **Regularization** (linear models).
- **“Fully”-Bayesian** methods combine both of these.
 - Average over all models, weighted by posterior (which includes regularizer).
 - Very powerful class of models we’ll cover in CPSC 540.
- Sometimes it’s **easier to define a likelihood than a loss function**.
 - We’ll do this for **multi-class classification**.

Multi-Label Classification

- We've been considering supervised learning with a **single label**:

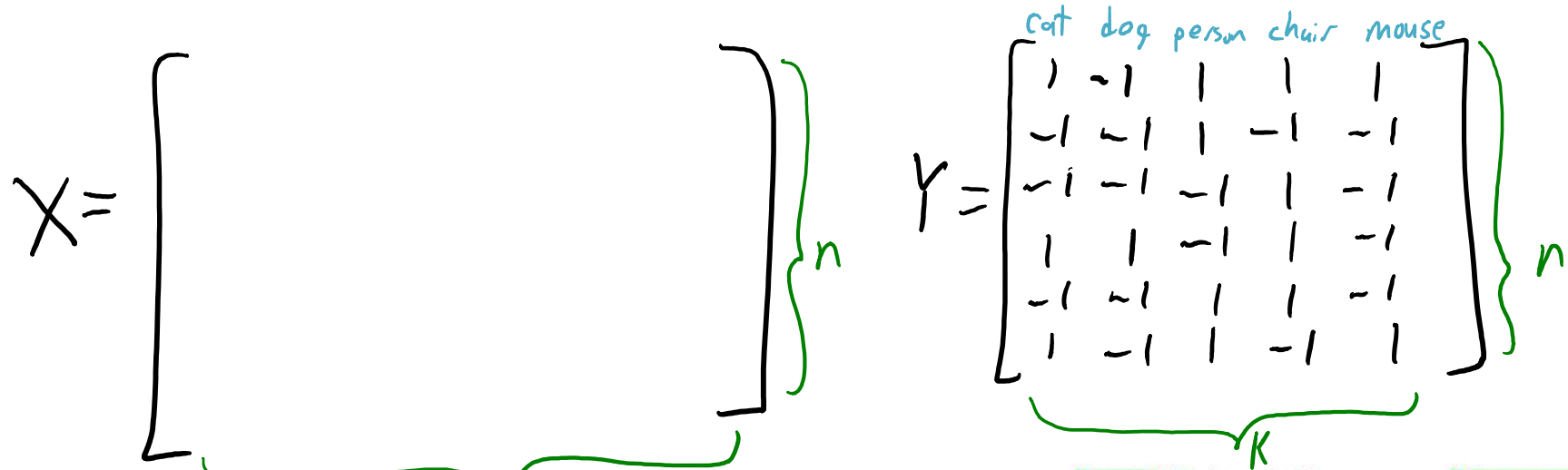
$$X = \begin{bmatrix} \\ \\ \\ \end{bmatrix} \quad y = \begin{bmatrix} 1 \\ -1 \\ -1 \\ -1 \end{bmatrix}$$

- E.g., is there a cat in this image or not?

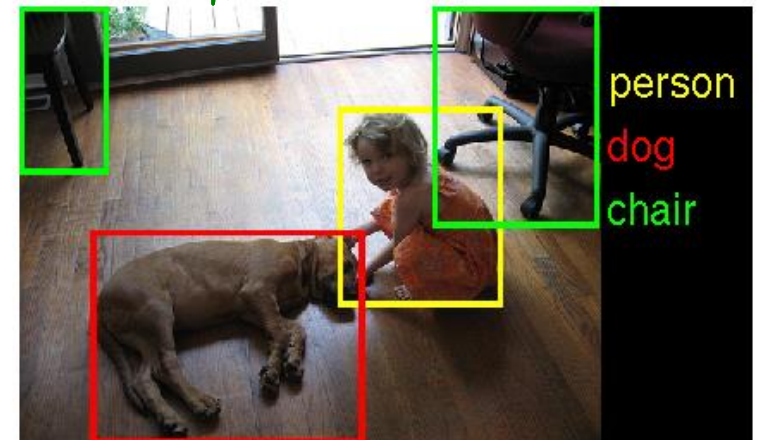


Multi-Label Classification

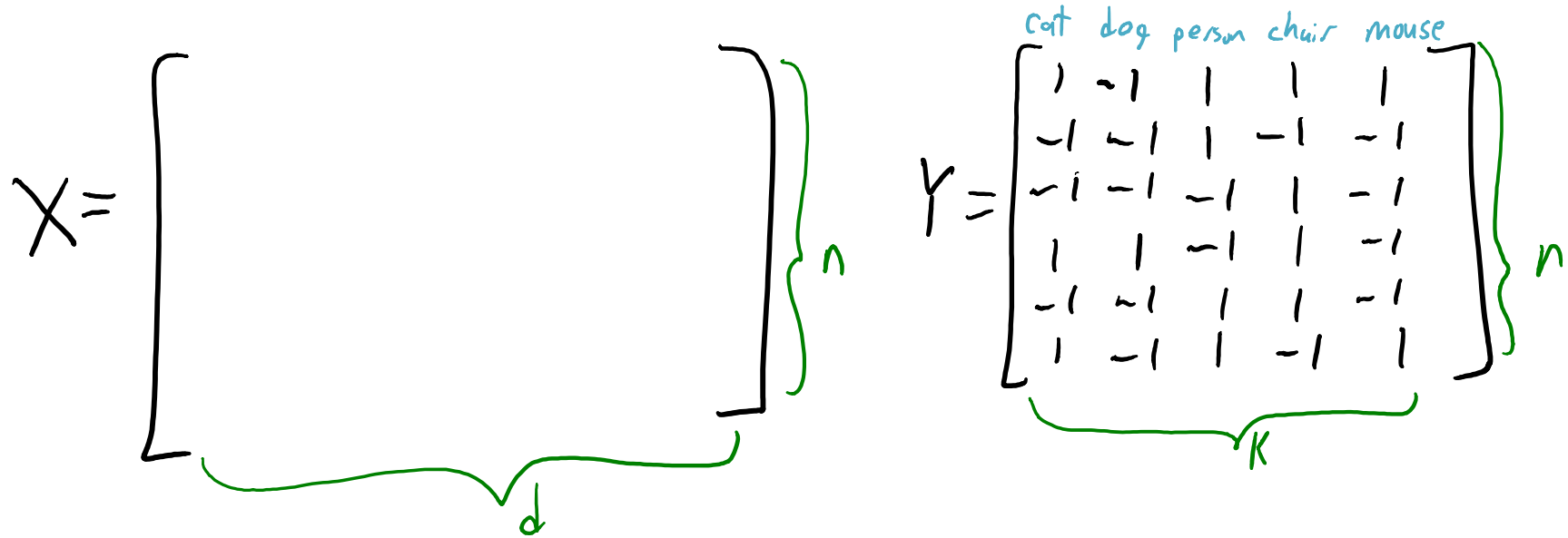
- In **multi-label classification** we want to predict '**k**' **binary labels**:



- E.g., which of the '**k**' objects are in this image?



Multi-Label Classification



- Approach 1:
 - Treat $\{1,-1,1,1,-1\}$ as the binary label 10110.
 - Problem is that with 'k' labels you have 2^k classes.
 - Only useful if 'k' is very small.

Multi-Label Classification

$$X = \left[\begin{array}{c} \\ \\ \\ \\ \\ \end{array} \right] \quad Y = \begin{array}{c} \text{cat} \quad \text{dog} \quad \text{person} \quad \text{chair} \quad \text{mouse} \\ \left[\begin{array}{ccccc} 1 & 1 & 1 & 1 & 1 \\ -1 & -1 & 1 & -1 & -1 \\ -1 & -1 & -1 & 1 & -1 \\ 1 & 1 & -1 & 1 & -1 \\ -1 & -1 & 1 & 1 & -1 \\ 1 & -1 & 1 & -1 & 1 \end{array} \right] \end{array}$$

Handwritten annotations: A green bracket on the right of X is labeled n . A green bracket below X is labeled d . A green bracket below the columns of Y is labeled K . A green bracket on the right of Y is labeled n .

- Approach 2:
 - Fit a binary classifier for each column 'c' of Y , using column as labels.
- If we use a linear model, each classifier has a weights w_c .
- Let's put the w_c together into a matrix 'W':

$$W = \left[\begin{array}{cccc} | & | & \dots & | \\ w_1 & w_2 & \dots & w_K \\ | & | & \dots & | \end{array} \right]$$

Handwritten annotations: A green bracket below the columns of W is labeled K . A green bracket on the right of W is labeled d .

Multi-Label Classification

$$X = \begin{bmatrix} \\ \\ \\ \\ \end{bmatrix} \quad Y = \begin{matrix} \text{cat} & \text{dog} & \text{person} & \text{chair} & \text{mouse} \\ \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ -1 & -1 & 1 & -1 & -1 \\ -1 & -1 & -1 & 1 & -1 \\ 1 & 1 & -1 & 1 & -1 \\ -1 & 1 & 1 & 1 & -1 \\ 1 & -1 & 1 & -1 & 1 \end{bmatrix} \end{matrix} \quad W = \begin{bmatrix} | & | & \dots & | \\ w_1 & w_2 & \dots & w_k \\ | & | & \dots & | \end{bmatrix}$$

Each column 'c' is a binary classifier for class 'c'

To predict label 'c' for example 'i' use $y_{ic} = \text{sign}(w_c^T x_i)$

To predict all labels y_i for example 'i' use $y_i = \text{sign}(W^T x_i)$

To predict all labels for all examples use $Y = \text{sign}(XW)$

- Fancier methods model correlations between y_i or between the w_c .

Multi-Class Classification

$$X = \begin{bmatrix} \\ \\ \\ \\ \\ \end{bmatrix} \quad Y = \begin{bmatrix} \text{cat} & \text{dog} & \text{person} & \text{chair} & \text{mouse} \\ 1 & -1 & -1 & -1 & -1 \\ -1 & -1 & 1 & -1 & -1 \\ -1 & -1 & -1 & 1 & -1 \\ -1 & -1 & -1 & 1 & -1 \\ -1 & -1 & -1 & 1 & -1 \\ 1 & -1 & -1 & -1 & -1 \end{bmatrix} \quad W = \begin{bmatrix} | & | & \dots & | \\ w_1 & w_2 & \dots & w_k \\ | & | & \dots & | \end{bmatrix}$$

Each column 'c' is a binary classifier for class 'c'

- **Multi-class classification:** special case where each y_i has 1 non-zero.

Multi-Class Classification and “One vs. All”

$$X = \begin{bmatrix} \\ \\ \\ \\ \end{bmatrix} \quad y = \begin{bmatrix} 1 \\ 3 \\ 4 \\ 4 \\ 1 \end{bmatrix}$$

$$W = \begin{bmatrix} | & | & \dots & | \\ w_1 & w_2 & \dots & w_k \\ | & | & \dots & | \end{bmatrix}$$

Each column 'c' is a binary classifier for class 'c'

- **Multi-class classification**: special case where each y_i has 1 non-zero.
 - Now we can code y_i as a discrete number $\{1, 2, 3, \dots\}$ giving class 'k'.
- **One vs. all** multi-class approach uses naïve multi-label approach:
 - Independently fit parameters ' w_c ' of a linear model for each class 'c'.
 - Each ' w_c ' tries to predict +1 for class 'c' and -1 for all others.

Multi-Class Classification and "One vs. All"

$$X = \begin{bmatrix} \\ \\ \\ \\ \end{bmatrix} \quad y = \begin{bmatrix} 1 \\ 3 \\ 4 \\ 4 \\ 1 \end{bmatrix}$$

$$W = \begin{bmatrix} | & | & \dots & | \\ w_1 & w_2 & \dots & w_k \\ | & | & \dots & | \end{bmatrix}$$

Each column 'c' is a binary classifier for class 'c'

- But prediction $W^T x_i$ might have multiple +1 values.
- To predict the "best" label, choose 'c' with largest value of $w_c^T x_i$.

$$y_i = \underset{c}{\operatorname{argmax}} \{ w_c^T x_i \}$$

Choose the 'c' that achieves the maximum value.

For logistic regression this choose class 'c' with the highest probability $p(y_{ic} = +1 | x_i, w)$

Multi-Class Classification and "One vs. All"

$$X = \begin{bmatrix} \\ \\ \\ \\ \end{bmatrix} \quad y = \begin{bmatrix} 1 \\ 3 \\ 4 \\ 4 \\ 1 \end{bmatrix}$$

$$W = \left[\begin{array}{c|c|c|c} | & | & \dots & | \\ w_1 & w_2 & & w_k \\ | & | & & | \end{array} \right] \Bigg\} \begin{array}{l} \text{Each column 'c' is a binary} \\ \text{classifier for class 'c'} \end{array}$$

$$y_i = \operatorname{argmax}_c \{ w_c^T x_i \}$$

Choose the 'c' that achieves the maximum value.

- But we **only trained w_c to get the correct sign of y_{ic}** :
 - We didn't train the w_c so that the largest $w_c^T x_i$ would be y_i .

Multinomial Logistic Regression

- Can we define a loss function so that largest $w_c^T x_i$ gives y_i ?
- In the **multi-label logistic regression** model we used:

$$p(y_i = +1 | x_i, W) = \frac{1}{1 + \exp(-w_c^T x_i)} = \frac{\exp(w_c^T x_i)}{\exp(w_c^T x_i) + 1} \propto \exp(w_c^T x_i)$$

- The **multinomial logistic regression** model uses the same idea:

$$p(y_i = c | x_i, W) \propto \exp(w_c^T x_i)$$

- But now need to **sum over 'k' classes** to get a valid probability:

$$p(y_i = c | x_i, W) = \frac{\exp(w_c^T x_i)}{\exp(w_1^T x_i) + \exp(w_2^T x_i) + \dots + \exp(w_k^T x_i)}$$

Multinomial Logistic Regression

- So **multinomial logistic regression** uses:

$$p(y_i | x_i, W) = \frac{\exp(w_{y_i}^T x_i)}{\sum_{c=1}^k \exp(w_c^T x_i)}$$

- Which is also known as the **softmax** function.
- Now that we have a probability, the MLE gives a loss function:

$$f(w) = - \sum_{i=1}^n \log(p(y_i | x_i, W))$$

$$= \sum_{i=1}^n \underbrace{-w_{y_i}^T x_i}_{\text{Tries to make this large for correct label.}} + \underbrace{\log\left(\sum_{c=1}^k \exp(w_c^T x_i)\right)}_{\text{Log-sum-exp which is smooth approximation to } \max_c \{\exp(w_c^T x_i)\} \text{ so we're trying to make the max small.}}$$

→ Convex so we can use gradient descent.

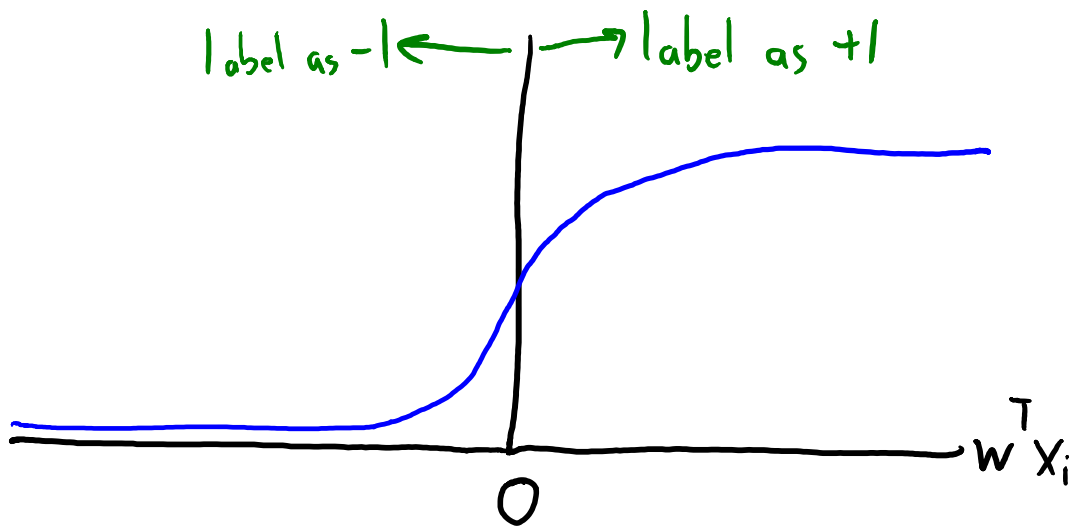
Losses for Other Discrete Labels

- MLE/MAP gives loss for classification with basic discrete labels:
 - Logistic regression for binary labels {"spam", "not spam"}.
 - Softmax regression for multi-class {"spam", "not spam", "important"}.
- But MLE/MAP lead to losses with other discrete labels:
 - Ordinal: {1 star, 2 stars, 3 stars, 4 stars, 5 stars}.
 - Counts: 602 'likes'.
- We can also use ratios of probabilities to define more losses:
 - Multi-class SVMs (similar to softmax, but generalizes hinge loss).
 - Ranking: $\text{Difficulty}(A3) > \text{Difficulty}(A4) > \text{Difficulty}(A2) > \text{Difficulty}(A1)$.

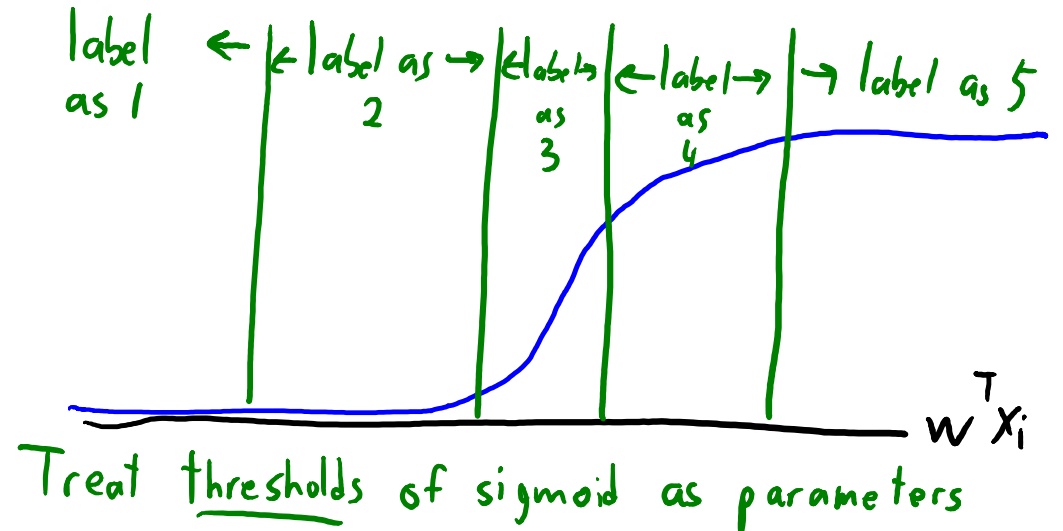
Ordinal Labels

- **Ordinal data**: categorical data where the **order matters**:
 - Rating hotels as {'1 star', '2 stars', '3 stars', '4 stars', '5 stars'}.
 - **Softmax would ignore order**.
- Can use '**ordinal logistic regression**'.

Logistic regression



Ordinal logistic regression



Count Labels

- **Count data**: predict the **number of times** something happens.
 - For example, $y_i = \text{“602”}$ Facebook likes.
- Softmax/ordinal **require finite number of possible labels**.
- We probably don't want separate parameter for '654' and '655'.
- **Poisson regression**: use probability from Poisson count distribution.
 - Many variations exist.

Summary

- **-log(probability)** lets us to define loss from any probability.
 - Special cases are least squares, least absolute error, and logistic regression.
- **MAP estimation** directly models $p(w \mid X, y)$.
 - Gives probabilistic interpretation to regularization.
- **Softmax loss** is natural generalization of logistic regression.
- **Discrete losses for weird scenarios** are possible using MLE/MAP:
 - Ordinal logistic regression, Poisson regression.
- Next time:
 - What ‘parts’ are your personality made of?

Bonus Slide: Multi-Class SVMs

Bonus Slide: Ranking