

# CPSC 340: Machine Learning and Data Mining

Multi-Class Regression

# Admin

- Midterm:
  - Today is the last day for grading concerns.
- Assignment 4:
  - Start now/soon!!!
- Tutorials next week
  - Topic is maximum likelihood and hw4 Q2.2

# Last Time: L1-Regularization

- We discussed L1-regularization:

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

- 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.

# Things I wish I said last time

- L1 regularization is **not** the same as the L1 loss
  - One is for sparsity in `w` and one is for robustness of fit, respectively
- The transition from discussing regularization to discussing maximum likelihood is not completely a non-sequitur
  - We'll see more of this today
  - There's an interpretation of regularization as a prior distribution on `w`

# Last Time: 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:

$$\begin{aligned} f(w) &= -\sum_{i=1}^n \log(p(y_i | x_i, w)) \\ &= -\sum_{i=1}^n \log\left(\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(w^T x_i - y_i)^2}{2\sigma^2}\right)\right) \\ &= -\sum_{i=1}^n \left[ \log\left(\frac{1}{\sqrt{2\pi\sigma^2}}\right) + \log\left(\exp\left(-\frac{(w^T x_i - y_i)^2}{2\sigma^2}\right)\right) \right] \end{aligned}$$

Constant in terms of  $w$

operations cancel

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

↓ 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:

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

- So least squares is MLE under Gaussian likelihood.
  - This is a big deal!

# Problem with Maximum Likelihood Estimation

- Maximum likelihood estimate maximizes:

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

- It's is 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 overfits 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

$$\begin{aligned} p(w | X, y) &= \frac{p(y | X, w) p(w | X)}{p(y | X)} \propto p(y | X, w) p(w | X) \\ &= p(y | X, w) p(w) \end{aligned}$$

↓ Assume ' $w$ ' does  
not depend on ' $X$ '!

"likelihood" "prior"

# 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  $p(h)$  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:

$$\begin{aligned} -\log(p(w | X, y)) &= -\log(p(y | X, w)) - \log(p(w)) + (\text{constant}) \\ &\quad \downarrow \text{IID data} \qquad \qquad \qquad \downarrow \text{prior over the } w_j \text{ are independent} \\ &= -\sum_{i=1}^n \log(p(y_i | x_i, w)) - \sum_{j=1}^d \log(p(w_j)) + (\text{constant}) \end{aligned}$$

# 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.

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

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

↑ for MAP estimation  
the constant changes  
solution.

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

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

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

# Summarizing...

- The choice of **likelihood** corresponds to the choice of **loss**
  - Our assumptions about how the  $y$ -values arise given the “truth”
- The choice of **prior** corresponds to the choice of **regularizer**
  - Our assumptions about ‘ $w$ ’ -- does not involve the data

# 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 covered 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**:

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

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

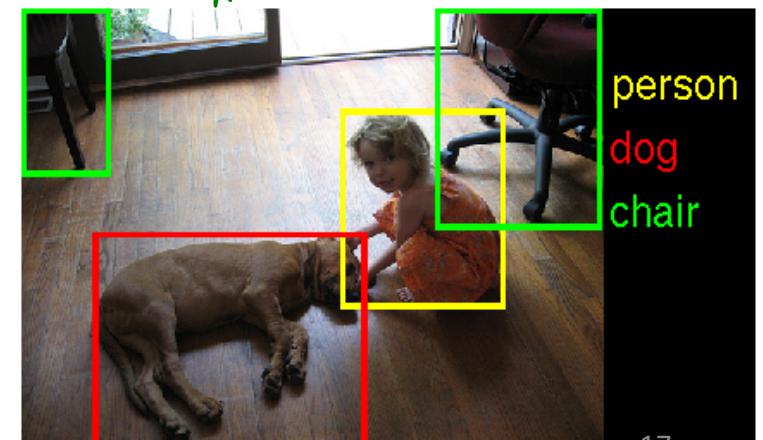


# Multi-Label Classification

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

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

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



# Multi-Label Classification

$$X = \begin{bmatrix} & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \end{bmatrix}^n$$

$\downarrow d$

$$Y = \begin{bmatrix} \text{cat} & \text{dog} & \text{person} & \text{chair} & \text{mouse} \\ 1 & -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & -1 & -1 \\ -1 & -1 & -1 & 1 & -1 \\ 1 & 1 & -1 & 1 & -1 \\ -1 & -1 & 1 & 1 & -1 \\ 1 & -1 & 1 & -1 & 1 \end{bmatrix}^n$$

$\downarrow K$

- 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

$$X = \begin{bmatrix} & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \end{bmatrix}^n$$
$$Y = \begin{bmatrix} \text{cat} & \text{dog} & \text{person} & \text{chair} & \text{mouse} \\ 1 & -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & -1 & -1 \\ -1 & -1 & -1 & 1 & -1 \\ 1 & 1 & -1 & 1 & -1 \\ -1 & -1 & 1 & 1 & -1 \\ 1 & -1 & 1 & -1 & 1 \end{bmatrix}^n$$

$d$        $n$        $K$

- 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’:  $W = \left[ \begin{array}{c|c|c|c|c} | & | & | & | & | \\ \hline w_1 & w_2 & \cdots & w_k \\ | & | & \cdots & | \end{array} \right] \{ c \}$

$$W = \left[ \begin{array}{c|c|c|c|c} | & | & | & | & | \\ \hline w_1 & w_2 & \cdots & w_k \\ | & | & \cdots & | \end{array} \right] \{ c \}$$

$K$

# Multi-Label Classification

$$X = \left[ \begin{array}{c} \\ \\ \\ \\ \end{array} \right]$$
$$Y = \left[ \begin{array}{ccccc} \text{cat} & \text{dog} & \text{person} & \text{chair} & \text{mouse} \\ \downarrow & \downarrow & \downarrow & \downarrow & \downarrow \\ 1 & -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & -1 & -1 \\ -1 & -1 & -1 & 1 & -1 \\ 1 & 1 & -1 & 1 & -1 \\ -1 & -1 & 1 & 1 & -1 \\ 1 & -1 & 1 & -1 & 1 \end{array} \right] \quad d$$
$$W = \left[ \begin{array}{cccc} | & | & \cdots & | \\ w_1 & w_2 & \cdots & w_K \\ | & | & \cdots & | \end{array} \right] \quad c$$

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  $w_c$ . 20

# Multi-Class Classification

$$X = \left[ \begin{array}{c} \\ \\ \\ \\ \end{array} \right] \quad Y = \left[ \begin{array}{ccccc} \text{cat} & \text{dog} & \text{person} & \text{chair} & \text{mouse} \\ 1 & -1 & -1 & -1 & -1 \\ -1 & -1 & 1 & -1 & -1 \\ -1 & -1 & -1 & 1 & -1 \\ -1 & -1 & -1 & 1 & -1 \\ -1 & -1 & -1 & -1 & 1 \end{array} \right]^d \quad W = \left[ \begin{array}{c} w_1 \\ w_2 \\ \vdots \\ w_K \end{array} \right]^c \quad \text{Each column } 'c' \text{ is a binary classifier for class } 'c'$$

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

# Multi-Class Classification and “One vs. All”

$$X = \left[ \begin{array}{c} \\ \\ \\ \\ \end{array} \right] \quad y = \left[ \begin{array}{c} 1 \\ 3 \\ 4 \\ 4 \\ 4 \\ 1 \end{array} \right]$$

$$W = \left[ \begin{array}{cccc} | & | & \cdots & | \\ w_1 & w_2 & \cdots & w_K \\ | & | & \cdots & | \end{array} \right]_C^K$$

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”

$$X = \left[ \begin{array}{c} \\ \\ \\ \\ \end{array} \right] \quad y = \left[ \begin{array}{c} 1 \\ 3 \\ 4 \\ 4 \\ 4 \\ 1 \end{array} \right]$$

$$W = \left[ \begin{array}{cccc} | & | & \cdots & | \\ w_1 & w_2 & \cdots & w_K \\ | & | & \cdots & | \end{array} \right] \}_{\text{c}}$$

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”

$$X = \left[ \begin{array}{c} \\ \\ \\ \\ \end{array} \right] \quad y = \left[ \begin{array}{c} 1 \\ 3 \\ 4 \\ 4 \\ 4 \\ 1 \end{array} \right]$$

$$W = \left[ \begin{array}{c|c|c|c} 1 & 1 & \cdots & 1 \\ \hline w_1 & w_2 & \cdots & w_K \end{array} \right] \}_{\text{c}}$$

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

$$y_i = \arg \max_c \{ w_c^T x_i \}$$

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:

$$p(y_{ic} = +1 | x_i, w) = \frac{1}{1 + \exp(-w_c^T x_i)} = \frac{\exp(w_c^T x_i)}{\exp(w_c^T x_i) + 1} \propto \exp(w_c^T x_i) \quad \text{and } p(y_i = -1 | x_i, w) \propto 1$$

- 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:

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

$$= \sum_{i=1}^n \left[ -w_{y_i}^T x_i + \log \left( \sum_{c=1}^k \exp(w_c^T x_i) \right) \right]$$

Tries to make  
this large for  
correct label.

Convex so we  
can use gradient  
descent.

Log-sum-exp which is smooth approximation  
to  $\max_c \{ \exp(w_c^T x_i) \}$  so we're trying to  
make the max small.

# Softmax Loss Function

- The softmax loss function:

$$f(w) = \sum_{i=1}^n \left[ -w_{y_i}^\top x_i + \log \left( \sum_{c=1}^k \exp(w_c^\top x_i) \right) \right]$$

$w_{y_i}$  are the weights for the true label.

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

These are the weights for label 'c'.

If  $y_i = 2$ ,  $w_{y_i}$  is  $w_2$ .

We have a column  $w_c$  giving the weights for class 'c'.

# 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'.
  - Survival rate: 60% of patients were still alive after 3 years.
- Define likelihood of labels, and NLL gives loss function.
- Later we'll use ratios of probabilities to define more losses:
  - Multi-class SVMs (similar to softmax, but generalizes hinge loss).
  - Ranking: Difficulty(A3) > Difficulty(A4) > Difficulty (A2) > Difficulty(A1).

# Summary

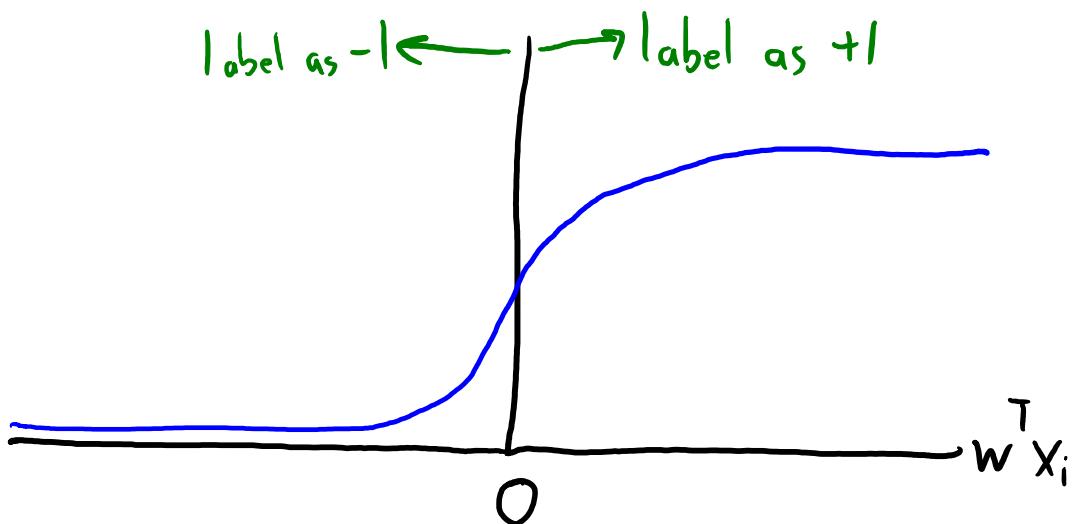
- **-log(probability)** lets us to define loss from any probability.
  - Special cases are least squares, least absolute error, and logistic regression.
- **MAP estimation** directly models  $p(w | X, y)$ .
  - Gives probabilistic interpretation to regularization.
- **Softmax loss** is natural generalization 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: Independence of X and w

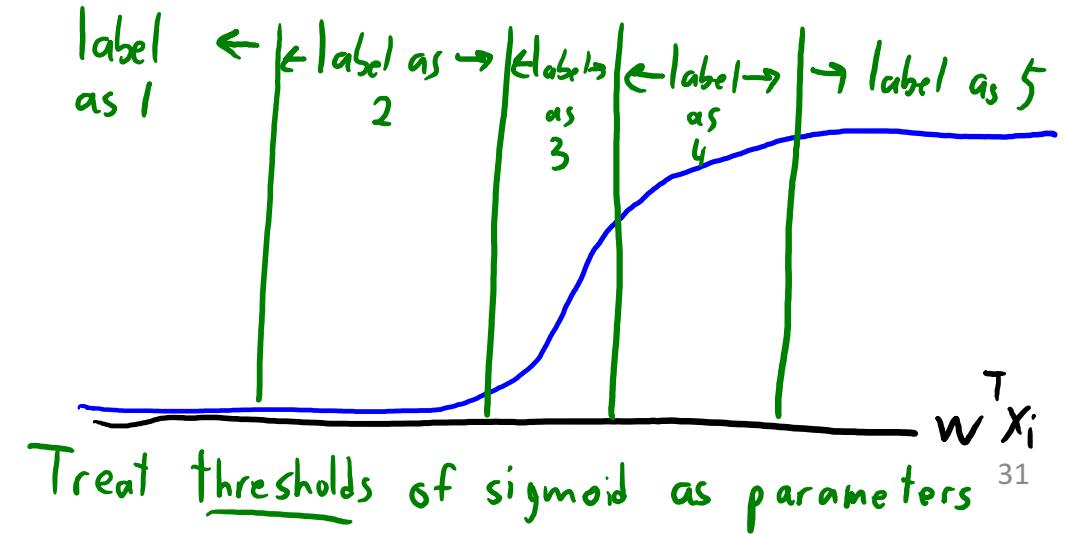
# Bonus Slide: Ordinal Labels

- **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



Treat thresholds of sigmoid as parameters

# Count Labels

- Count data: predict the number of times something happens.
  - For example,  $y_i = "602"$  Facebook likes.
- Softmax requires 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.