Linear Regression (Contd), Linear Classification

CS771: Introduction to Machine Learning Piyush Rai

Gradient Descent for Linear/Ridge Regression

Just use the GD algorithm with the gradient expressions we derived

 $w^{(t+1)} = w^{(t)} - n a^{(t)}$

Iterative updates for linear regression will be of the form

Also, we usually work with average gradient so the gradient term is divided by N

Note the form of each term in the gradient expression update: Amount of current w's error on the n^{th} training example multiplied by the input x_n

Unlike the closed form solution

$$(\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}}\mathbf{y}$$
 of least squares
regression, here we have iterative
updates but do not require the
expensive matrix inversion of the
 $D \times D$ matrix $\mathbf{X}^{\mathsf{T}}\mathbf{X}$

$$= \mathbf{w}^{(t)} + \eta_t \sum_{n=1}^{N} \left(y_n - \mathbf{w}^{(t)}^{\mathsf{T}}\mathbf{x}_n \right) \mathbf{x}_n$$

- Similar updates for ridge regression as well (with the gradient expression being slightly different; left as an exercise)
- More on iterative optimization methods later



 ℓ_2 regularization and "Smoothness"

The regularized objective we minimized is

$$L_{reg}(\boldsymbol{w}) = \sum_{n=1}^{N} (y_n - \boldsymbol{w}^{\mathsf{T}} \boldsymbol{x}_n)^2 + \lambda \boldsymbol{w}^{\mathsf{T}} \boldsymbol{w}$$

- Minimizing $L_{reg}(w)$ w.r.t. w gives a solution for w that
 - Keeps the training error small
 - Has a small ℓ_2 squared norm $\mathbf{w}^{\mathsf{T}}\mathbf{w} = \sum_{d=1}^{D} w_d^2$

Remember – in general, weights with large magnitude are bad since they can cause overfitting on training data and may not work well on test data



Not a "smooth" model since its test data predictions may change drastically even with small changes in some feature's value

Small entries in \boldsymbol{w} are good since they lead to "smooth" models



CS771: Intro to ML

Good because, consequently, the individual entries of the weight vector \boldsymbol{w} are also prevented from becoming too large

Other Ways to Control Overfitting

• Use a regularizer R(w) defined by other norms, e.g.,

Note that optimizing loss functions with such regularizers is usually harder than ridge reg. but several advanced techniques exist (we will see some of those later)

Use them if you have a very large number of features but regularizers can help in

influence prediction

- Use non-regularization based approaches
 - Early-stopping (stopping training just when we have a decent val. set accuracy)
 - Dropout (in each iteration, don't update some of the weights)
 - Injecting noise in the inputs

All of these are very popular ways to control overfitting in deep learning models. More on these later when we talk about deep learning



Linear Regression as Solving System of Linear Eqs

- The form of the lin. reg. model $y \approx Xw$ is akin to a system of linear equation
- Assuming N training examples with D features each, we have

First training example: $y_1 = x_{11}w_1 + x_{12}w_2 + \dots + x_{1D}w_D$ Second training example: $y_2 = x_{21}w_1 + x_{22}w_2 + \dots + x_{2D}w_D$ Note: Here x_{nd} denotes the d^{th} feature of the n^{th} training example

```
N equations and D unknowns
here (w_1, w_2, ..., w_D)
```

Now solve this!

 $y_N = x_{N1}w_1 + x_{N2}w_2 + \dots + x_{ND}w_D$ N-th training example:

- Usually we will either have N > D or N < D
 - Thus we have an underdetermined (N < D) or overdetermined (N > D) system
 - Methods to solve over/underdetermined systems can be used for lin-reg as well
 - Many of these methods don't require expensive matrix inversion

Solving lin-reg as system of lin eq. $w = (X^{\mathsf{T}}X)^{-1} X^{\mathsf{T}}y \qquad \longrightarrow \qquad Aw = b$ where $A = X^{\mathsf{T}}X$, and $b = X^{\mathsf{T}}y$

System of lin. Eqns with D equations and D unknowns

The bias term

• Linear models usually also have a bias term b in addition to the weights



Can append a constant feature "1" for each input and again rewrite as $y = \widetilde{w}^{\mathsf{T}}\widetilde{x}$ where now both $\widetilde{x} = [1, x]$ and $\widetilde{w} = [b, w]$ are in \mathbb{R}^{D+1}

We will assume the same and omit the explicit bias for simplicity of notation

D + 1



Evaluation Measures for Regression Models Prediction

- Plotting the prediction \hat{y}_n vs truth y_n for the validation/test set
- Mean Squared Error (MSE) and Mean Absolute Error (MAE) on val./test set

$$MSE = \frac{1}{N} \sum_{n=1}^{N} (y_n - \hat{y}_n)^2 \qquad MAE = \frac{1}{N} \sum_{n=1}^{N} |y_n - \hat{y}_n|^2$$

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

- RMSE (Root Mean Squared Error) $\triangleq \sqrt{MSE}$
- Coefficient of determination or R^2



Linear Models for Classification



Linear Models for Classification

- A linear model $y = w^{T}x$ can also be used in classification
- For binary classification, can treat $w^{T}x_{n}$ as the "score" of input x_{n} and either



• Note: In LR, if we assume the label y_n as -1/+1 (not 0/1) then we can write $p(y_n|w, x_n) = \frac{1}{1 + \exp(-y_n w^{\mathsf{T}} x_n)} = \sigma(y_n w^{\mathsf{T}} x_n)$

Linear Models: The Decision Boundary

• Decision boundary is where the score $w^{\mathsf{T}} x_n$ changes its sign



Therefore, both views are equivalent

- Decision boundary is where both classes have equal probability for the input x_n
- For logistic reg, at decision boundary $p(y_n = 1 | \boldsymbol{w}, \boldsymbol{x}_n) = p(y_n = 0 | \boldsymbol{w}, \boldsymbol{x}_n)$ $\frac{\exp(\boldsymbol{w}^{\mathsf{T}} \boldsymbol{x}_n)}{1 + \exp(\boldsymbol{w}^{\mathsf{T}} \boldsymbol{x}_n)} = \frac{1}{1 + \exp(\boldsymbol{w}^{\mathsf{T}} \boldsymbol{x}_n)}$ $\exp(\boldsymbol{w}^{\mathsf{T}} \boldsymbol{x}_n) = 1$ $\boldsymbol{w}^{\mathsf{T}} \boldsymbol{x}_n = 0$



Linear Models for (Multi-class) Classification

• If there are K > 2 classes, we use K weight vectors $\{w_i\}_{i=1}^K$ to define the model

The prediction rule is as follows

$$y_n = \operatorname{argmax}_{i \in \{1,2,\dots,K\}} \boldsymbol{w}_i^{\mathsf{T}} \boldsymbol{x}_n$$

- Can think of $\boldsymbol{w}_i^{\mathsf{T}} \boldsymbol{x}_n$ as the score/similarity of the input w.r.t. the i^{th} class
- Can also use these scores to compute probability of belonging to each class

$$\mu_{n,i} = p(y_n = i | W, x_n) = \frac{\exp(w_i^{\mathsf{T}} x_n)}{\sum_{j=1}^{K} \exp(w_j^{\mathsf{T}} x)} \text{ Multi-class extension} \text{ of logistic regression} \qquad \mu_{n,i} \qquad \mu_{n,i}$$

Linear Classification: Interpreting weight vectors

Recall that multi-class classification prediction rule is

$$y_n = \operatorname{argmax}_{i \in \{1,2,\dots,K\}} \boldsymbol{w}_i^{\mathsf{T}} \boldsymbol{x}_n$$

- Can think of $w_i^T x_n$ as the score of the input for the i^{th} class (or similarity of x_n with w_i)
- Once learned (we will see the methods later), these K weight vectors (one for each class) can sometimes have nice interpretations, especially when the inputs are images



Loss Functions for Classification

- Assume true label to be $y_n \in \{0,1\}$ and the score of a linear model to be $w^{\top}x_n$
- One possibility is to use squared loss just like we used in regression

$$l(y_n, \boldsymbol{w}^{\mathsf{T}}\boldsymbol{x}_n) = (y_n - \boldsymbol{w}^{\mathsf{T}}\boldsymbol{x}_n)^2$$

- Will be easy to optimize (same solution as the regression case)
- Can also consider other loss functions used in regression
 - Basically, pretend that the binary label is actually a continuous value and treat the problem as regression where the output can only be one of two possible values
- However, regression loss functions aren't ideal since y_n is discrete (binary/categorical)
- Using the score $w^T x_n$ or the probability $\mu_n = \sigma(w^T x_n)$ of belonging to the positive class, we have specialized loss function for binary classification

Loss Functions for Classification: Cross-Entropy

- Cross-entropy (CE) is a popular loss function for binary classification. Used in logistic reg.
- Assuming true $y_n \in \{0,1\}$ and $\mu_n = \sigma(\mathbf{w}^\top \mathbf{x}_n)$ as predicted prob of $y_n = 1$, CE loss is Cross-Entropy Loss for Binary Classification

$$L(\mathbf{w}) = -\left[\sum_{n=1}^{N} y_n \log \mu_n + (1 - y_n) \log(1 - \mu_n)\right]$$

Very large loss if y_n is 1 and μ_n close to 0, or y_n is 0and μ_n close to 1

later

For multi-class classification, the CE loss is defined as

$$L(W) = -\sum_{n=1}^{N} \sum_{i=1}^{K} y_{n,i} \log \mu_{n,i}$$

CE loss is also convex in **w**
(can prove easily using
definition of convexity; will see
later). Therefore unique solution
is obtained when we minimize it

Note: Sometimes we divide the loss function (not just CE but others too like squared loss) by the number of training examples N (doesn't make a difference to the solution; just a scaling factor. All relevant quantities, such as gradients will also get divided by N

0.2

04

Predicted Probability

LOSS

0.0

True Label = 0 True Label = 1

0.6

0.8

10



14

Cross-Entropy Loss: The Gradient

The expression for the gradient of binary cross-entropy loss

$$g = \nabla_{w} L(w) = -\sum_{n=1}^{N} (y_n - \mu_n) x_n$$
Using this, we can now do
gradient descent to learn the
optimal w for logistic regression:
$$w^{(t+1)} = w^{(t)} - \eta_t g^{(t)}$$
Note the form of each term in the gradient expression:
Amount of current w's error in predicting the label of
the nth training example multiplied by the input x_n

Note the u is a

CS771: Intro to ML

The expression for the gradient of multi-class cross-entropy loss

Linear Models for Classification

- A linear model $y = w^{T}x$ can also be used in classification
- For binary classification, can treat $w^{\mathsf{T}}x_n$ as the "score" of input x_n and either



Some Other Loss Functions for Binary Classification

• Assume true label as y_n and prediction as $\hat{y}_n = \operatorname{sign}[\boldsymbol{w}^{\mathsf{T}}\boldsymbol{x}_n]$

The zero-one loss is the most natural loss function for classification



Since zero-one loss is hard to minimize, we use some surrogate loss function

- Popular examples: Cross-entropy (same as logistic loss), hinge loss, etc
- Note: Ideally, surrogate loss (approximation of zero-one) must be an <u>upper bound</u> (must be larger than the O-1 loss for all values of $y_n w^T x_n$) since our goal is minimization.



Evaluation Measures for Binary Classification

Average classification error or average accuracy (on val./test data)

$$err(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^{N} \mathbb{I}[y_n \neq \hat{y}_n] \qquad acc(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^{N} \mathbb{I}[y_n = \hat{y}_n]$$

The cross-entropy loss itself (on val./test data)

- Precision, Recall, and F1 score (preferred if labels are imbalanced)
 - Precision (P): Of positive predictions by the model, what fraction is true positive
 - Recall (R): Of all true positive examples, what fraction the model predicted as positive
 - F1 score: Harmonic mean of P and R
- Confusion matrix is also a helpful measure



Various other metrics such as error/accuracy, P, R, F1, etc. can be readily calculated from the confusion matrix



Evaluation Measures for Multi-class Classification

Average classification error or average accuracy (on val./test data)

$$err(w) = \frac{1}{N} \sum_{n=1}^{N} \mathbb{I}[y_n \neq \hat{y}_n] \qquad acc(w) = \frac{1}{N} \sum_{n=1}^{N} \mathbb{I}[y_n = \hat{y}_n]$$

= Top-k accuracy

$$y_n \text{ is the true label, } \hat{s}_n \text{ is the set of top-k predicted classes for } x_n \text{ (based on the predicted probabilities/scores of the various classes)}$$

Top - k Accuracy =
$$\frac{1}{N} \sum_{n=1}^{N} \text{is_correct_top_k}[y_n, \hat{S}_n]$$

- The cross-entropy loss itself (on val./test data)
- Class-wise Precision, Recall, and F1 score (preferred if labels are imbalanced)
- Confusion matrix





 $\mathbb{I}[y_n = \hat{y}_n]$

Coming up next

Optimization techniques for machine learning



21