

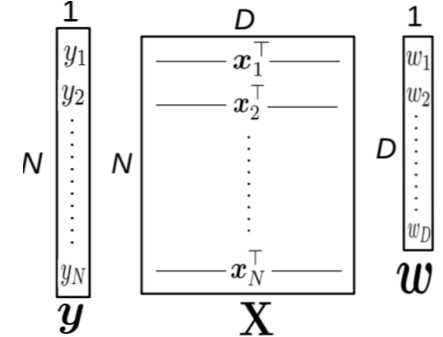
Linear Regression

CS771: Introduction to Machine Learning

Piyush Rai

Linear Regression

- Given: Training data with N input-output pairs $\{(\mathbf{x}_n, y_n)\}_{n=1}^N$, $\mathbf{x}_n \in \mathbb{R}^D$, $y_n \in \mathbb{R}$
- Goal: Learn a model to predict the output for new test inputs



- Assume the function that approximates the I/O relationship to be a linear model

$$y_n \approx f(\mathbf{x}_n) = \mathbf{w}^\top \mathbf{x}_n \quad (n = 1, 2, \dots, N)$$

Can also write all of them compactly using matrix-vector notation as $\mathbf{y} \approx \mathbf{X}\mathbf{w}$

- Let's write the total error or "loss" of this model over the training data as

Goal of learning is to find the \mathbf{w} that minimizes this loss + does well on test data

$$L(\mathbf{w}) = \sum_{n=1}^N \ell(y_n, \mathbf{w}^\top \mathbf{x}_n)$$

Unlike models like KNN and DT, here we have an explicit problem-specific objective (loss function) that we wish to optimize for

$\ell(y_n, \mathbf{w}^\top \mathbf{x}_n)$ measures the prediction error or "loss" or "deviation" of the model on a single training input (\mathbf{x}_n, y_n)

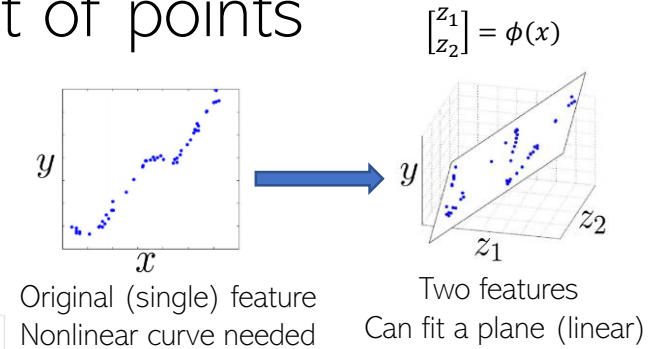
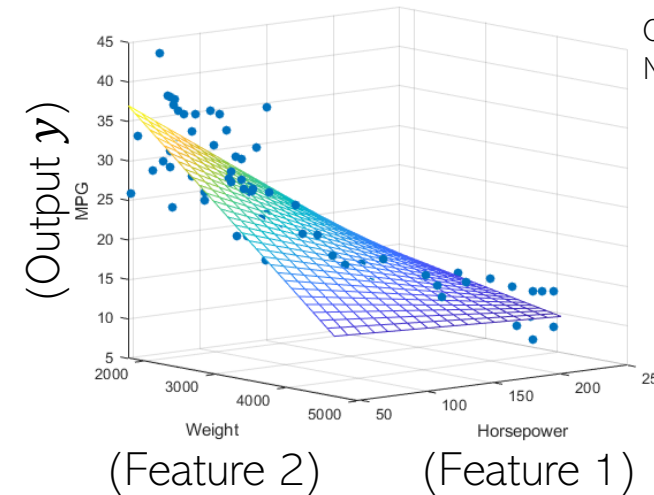
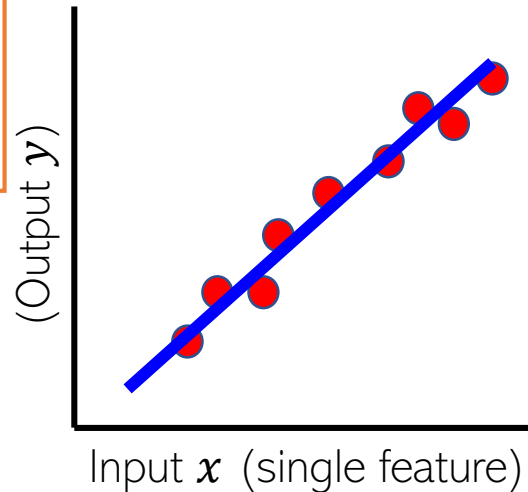


Linear Regression: Pictorially

- Linear regression is like fitting a line or (hyper)plane to a set of points

What if a line/plane doesn't model the input-output relationship very well, e.g., if their relationship is better modeled by a nonlinear curve or curved surface?

Do linear models become useless in such cases?



No. We can even fit a curve using a linear model after suitably transforming the inputs

$$y \approx \mathbf{w}^T \phi(x)$$

The transformation $\phi(\cdot)$ can be predefined or learned (e.g., using [kernel methods](#) or a deep neural network based feature extractor). More on this later

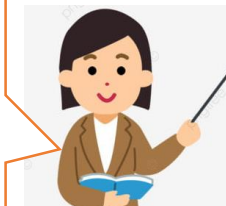
- The line/plane must also predict outputs the unseen (test) inputs well

Loss Functions for Regression

- Many possible loss functions for regression problems

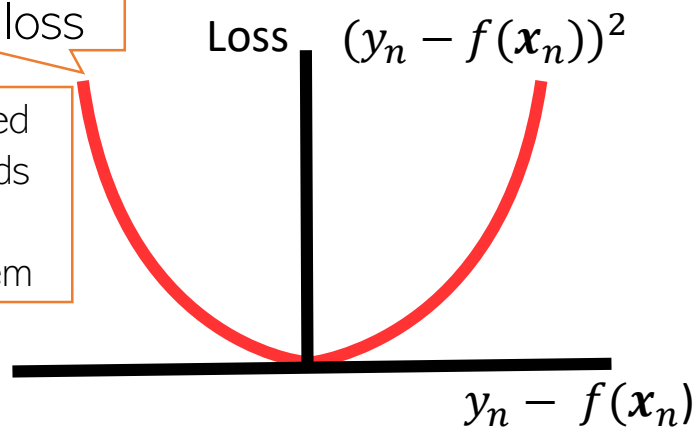
Choice of loss function usually depends on the nature of the data. Also, some loss functions result in easier optimization problem than others

4



Squared loss

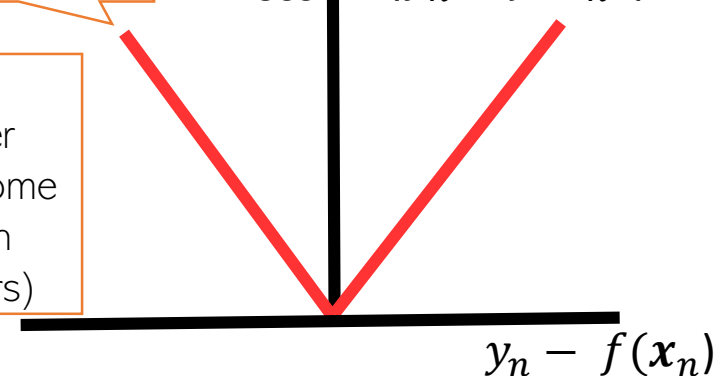
Very commonly used for regression. Leads to an easy-to-solve optimization problem



Absolute loss

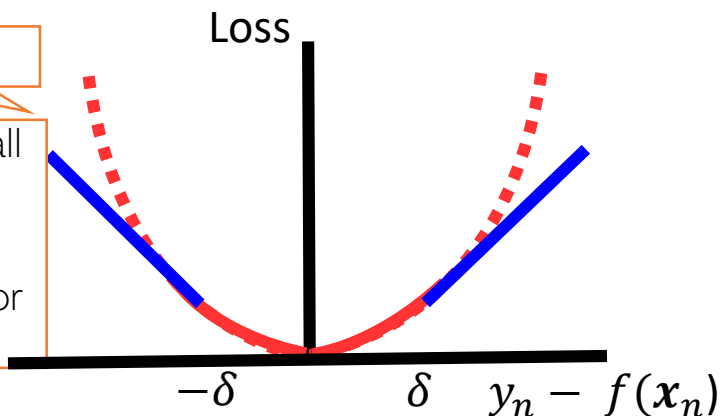
Grows more slowly than squared loss. Thus better suited when data has some outliers (inputs on which model makes large errors)

Loss $|y_n - f(x_n)|$



Huber loss

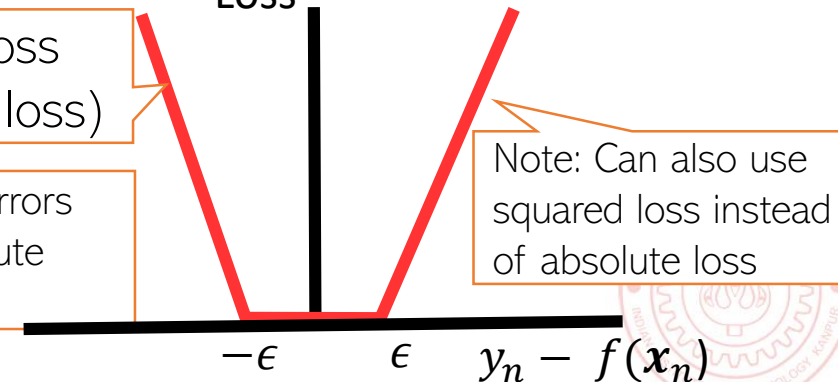
Squared loss for small errors (say up to δ); absolute loss for larger errors. Good for data with outliers



ϵ -insensitive loss (a.k.a. Vapnik loss)

Zero loss for small errors (say up to ϵ); absolute loss for larger errors

Loss $|y_n - f(x_n)| - \epsilon$



Note: Can also use squared loss instead of absolute loss

Linear Regression with Squared Loss

- In this case, the loss func will be

In matrix-vector notation, can write it compactly as $\|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 = (\mathbf{y} - \mathbf{X}\mathbf{w})^\top (\mathbf{y} - \mathbf{X}\mathbf{w})$

$$L(\mathbf{w}) = \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2$$

- Let us find the \mathbf{w} that optimizes (minimizes) the above squared loss

- We need calculus and optimization to do this!

The “least squares” (LS) problem
Gauss-Legendre, 18th century)

- The LS problem can be solved easily and has a closed form solution

$$\mathbf{w}_{LS} = \arg \min_{\mathbf{w}} L(\mathbf{w}) = \arg \min_{\mathbf{w}} \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2$$

Closed form solutions to ML problems are rare.

$$\mathbf{w}_{LS} = \left(\sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top \right)^{-1} \left(\sum_{n=1}^N y_n \mathbf{x}_n \right) = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

$D \times D$ matrix inversion – can be expensive.
Ways to handle this. Will see later



Proof: A bit of calculus/optim. (more on this later) ⁶

- We wanted to find the minima of $L(\mathbf{w}) = \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2$
- Let us apply basic rule of calculus: Take first derivative of $L(\mathbf{w})$ and set to zero

$$\frac{\partial L(\mathbf{w})}{\partial \mathbf{w}} = \frac{\partial}{\partial \mathbf{w}} \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2 = \sum_{n=1}^N 2(y_n - \mathbf{w}^\top \mathbf{x}_n) \frac{\partial}{\partial \mathbf{w}} (y_n - \mathbf{w}^\top \mathbf{x}_n) = 0$$

Chain rule of calculus

Partial derivative of dot product w.r.t each element of \mathbf{w}

Result of this derivative is \mathbf{x}_n - same size as \mathbf{w}

- Using the fact $\frac{\partial}{\partial \mathbf{w}} \mathbf{w}^\top \mathbf{x}_n = \mathbf{x}_n$, we get $\sum_{n=1}^N 2(y_n - \mathbf{w}^\top \mathbf{x}_n) \mathbf{x}_n = 0$
- To separate \mathbf{w} to get a solution, we write the above as

$$\sum_{n=1}^N 2\mathbf{x}_n (y_n - \mathbf{x}_n^\top \mathbf{w}) = 0 \quad \longrightarrow \quad \sum_{n=1}^N y_n \mathbf{x}_n - \mathbf{x}_n \mathbf{x}_n^\top \mathbf{w} = 0$$

$$\mathbf{w}_{LS} = \left(\sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top \right)^{-1} \left(\sum_{n=1}^N y_n \mathbf{x}_n \right) = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$



Problem(s) with the Solution!

- We minimized the objective $L(\mathbf{w}) = \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2$ w.r.t. \mathbf{w} and got

$$\mathbf{w}_{LS} = (\sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top)^{-1} (\sum_{n=1}^N y_n \mathbf{x}_n) = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

- Problem: The matrix $\mathbf{X}^\top \mathbf{X}$ may not be invertible
 - This may lead to non-unique solutions for \mathbf{w}_{opt}
- Problem: Overfitting since we only minimized loss defined on training data
 - Weights $\mathbf{w} = [w_1, w_2, \dots, w_D]$ may become arbitrarily large to fit training data perfectly
 - Such weights may perform poorly on the test data however
- One Solution: Minimize a **regularized objective** $L(\mathbf{w}) + \lambda R(\mathbf{w})$
 - The reg. will prevent the elements of \mathbf{w} from becoming too large
 - Reason: Now we are minimizing **training error** + **magnitude of vector \mathbf{w}**

$R(\mathbf{w})$ is called the **Regularizer** and measures the “magnitude” of \mathbf{w}

$\lambda \geq 0$ is the **reg. hyperparam.**
Controls how much we wish to regularize (needs to be tuned via cross-validation)

Regularized Least Squares (a.k.a. Ridge Regression)⁸

- Recall that the regularized objective is of the form $L_{reg}(\mathbf{w}) = L(\mathbf{w}) + \lambda R(\mathbf{w})$
- One possible/popular regularizer: the squared Euclidean (ℓ_2 squared) norm of \mathbf{w}

$$R(\mathbf{w}) = \|\mathbf{w}\|_2^2 = \mathbf{w}^\top \mathbf{w}$$

- With this regularizer, we have the regularized least squares problem as

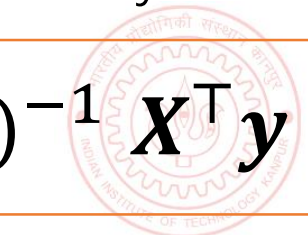
$$\mathbf{w}_{ridge} = \arg \min_{\mathbf{w}} L(\mathbf{w}) + \lambda R(\mathbf{w})$$

$$= \arg \min_{\mathbf{w}} \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2 + \lambda \mathbf{w}^\top \mathbf{w}$$

Look at the form of the solution. We are adding a small value λ to the diagonals of the $D \times D$ matrix $\mathbf{X}^\top \mathbf{X}$ (like adding a ridge/mountain to some land)

- Proceeding just like the LS case, we can find the optimal \mathbf{w} which is given by

$$\mathbf{w}_{ridge} = (\sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top + \lambda \mathbf{I}_D)^{-1} (\sum_{n=1}^N y_n \mathbf{x}_n) = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_D)^{-1} \mathbf{X}^\top \mathbf{y}$$



A closer look at ℓ_2 regularization

9

- The regularized objective we minimized is

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

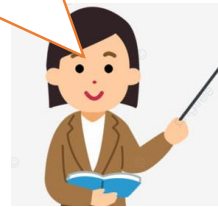
- Minimizing $L_{reg}(\mathbf{w})$ w.r.t. \mathbf{w} gives a solution for \mathbf{w} that

- Keeps the training error small
- Has a small ℓ_2 squared norm $\mathbf{w}^T \mathbf{w} = \sum_{d=1}^D w_d^2$

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

- Small entries in \mathbf{w} are good since they lead to “smooth” models

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

$\mathbf{x}_n =$	1.2	0.5	2.4	0.3	0.8	0.1	0.9	2.1
$\mathbf{x}_m =$	1.2	0.5	2.4	0.3	$0.8 + \epsilon$	0.1	0.9	2.1

Exact same feature vectors only differing in just one feature by a small amount

$$y_n = 0.8$$

$$y_m = 100$$

Very different outputs though (maybe one of these two training ex. is an outlier)

A typical \mathbf{w} learned without ℓ_2 reg.

3.2	1.8	1.3	2.1	10000	2.5	3.1	0.1
-----	-----	-----	-----	-------	-----	-----	-----

Just to fit the training data where one of the inputs was possibly an outlier, this weight became too big. Such a weight vector will possibly do poorly on normal test inputs

Other Ways to Control Overfitting

- Use a regularizer $R(\mathbf{w})$ defined by other norms, e.g.,

ℓ_1 norm regularizer

$$\|\mathbf{w}\|_1 = \sum_{d=1}^D |w_d|$$

$$\|\mathbf{w}\|_0 = \#\text{nnz}(\mathbf{w})$$

ℓ_0 norm regularizer (counts number of nonzeros in \mathbf{w})

When should I use these regularizers instead of the ℓ_2 regularizer?

Automatic feature selection? Wow, cool!!! But how exactly?

Use them if you have a very large number of features but many irrelevant features. These regularizers can help in [automatic feature selection](#)

Using such regularizers gives a [sparse](#) weight vector \mathbf{w} as solution

sparse means many entries in \mathbf{w} will be zero or near zero. Thus those features will be considered irrelevant by the model and will not 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

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)



Linear Regression as Solving System of Linear Eqs ¹¹

- The form of the lin. reg. model $\mathbf{y} \approx \mathbf{X}\mathbf{w}$ 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: $y_2 = x_{21}w_1 + x_{22}w_2 + \dots + x_{2D}w_D$

⋮

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

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, \dots, w_D)

- However, in regression, we rarely have $N = D$ but rather $N > D$ 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.

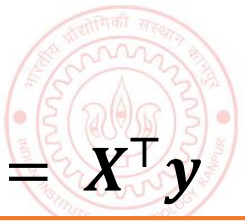
$$\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$



$$\mathbf{A}\mathbf{w} = \mathbf{b} \text{ where } \mathbf{A} = \mathbf{X}^T \mathbf{X}, \text{ and } \mathbf{b} = \mathbf{X}^T \mathbf{y}$$

Now solve this!

System of lin. Eqns with D equations and D unknowns



Next Lecture

- Solving linear regression using iterative optimization methods
 - Faster and don't require matrix inversion
- Brief intro to optimization techniques

