

# Linear Models and Learning via Optimization

CS771: Introduction to Machine Learning

Piyush Rai

# Announcements

- Quiz 1 postponed by a week (now on Aug 23)
  - Venue TBD, timing: 7pm, duration: 45 minutes
  - Syllabus will be everything up to Aug 21 class
- Homework 1 released by end of this week



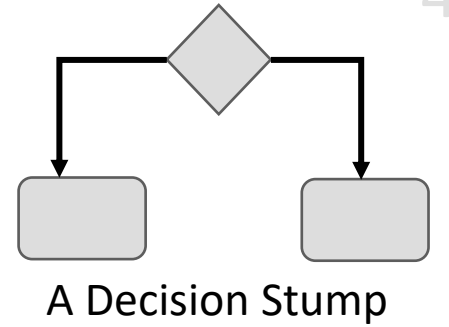
# Wrapping Up Decision Trees..



# Avoiding Overfitting in DTs

4

- Desired: a DT that is not too big in size, yet fits the training data reasonably
- Note: An example of a very simple DT is “decision-stump”
  - A decision-stump only tests the value of a single feature (or a simple rule)
  - Not very powerful in itself but often used in a large ensemble of decision stumps
- Some ways to keep a DT simple enough:
  - Control its complexity while building the tree (stopping early)
  - Prune after building the tree (post-pruning)
- Criteria for judging which nodes could potentially be pruned (from already built complex DT)
  - Use a validation set (separate from the training set)
    - Prune each possible node that doesn't hurt the accuracy on the validation set
    - Greedily remove the node that improves the validation accuracy the most
    - Stop when the validation set accuracy starts worsening



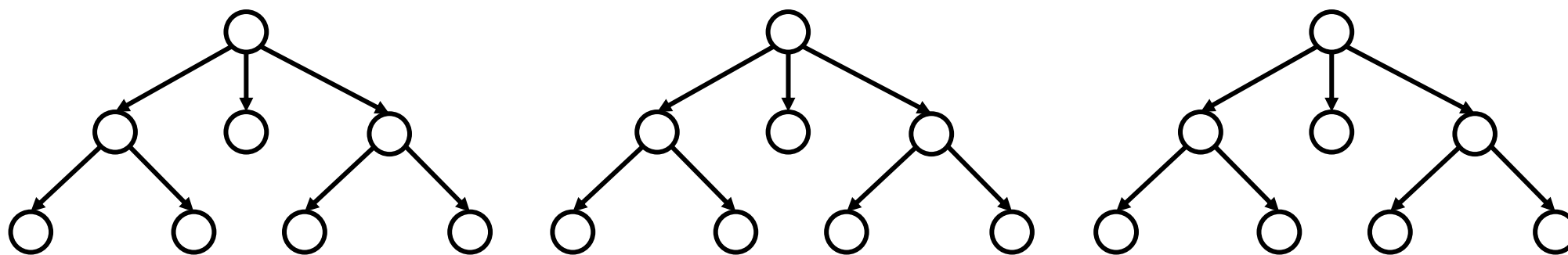
# Ensemble of Trees

- Ensemble is a collection of models. Popular in ML
- Each model makes a prediction. Take their majority as the final prediction
- Ensemble of trees is a collection of simple DTs
  - Often preferred as compared to a single massive, complicated tree
- A popular example: **Random Forest (RF)**

All trees can be trained in parallel

Each tree is trained on a subset of the training inputs/features

An RF with 3 simple trees. The majority prediction will be the final prediction



- **XGBoost** is another popular ensemble of trees
  - Based on the idea of “**boosting**” (will study boosting later) simple trees
  - Sequentially trains a set of trees with each correcting errors of previous ones



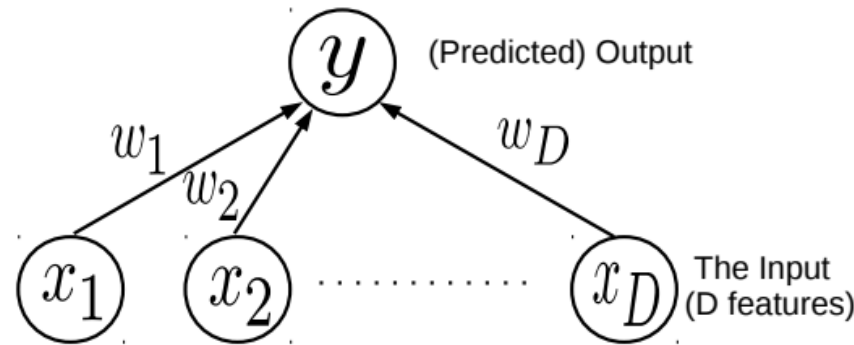
# Linear Models and Learning via Optimization



# Linear Models

- Suppose we want to learn to map inputs  $\mathbf{x} \in \mathbb{R}^D$  to real-valued outputs  $y \in \mathbb{R}$
- Linear model: Assume output to be a linear **weighted combination** of the  $D$  input features

$$y = \sum_{d=1}^D w_d x_d = \mathbf{w}^\top \mathbf{x}$$



This defines a linear model with  $D$  parameters given by a “weight vector”  $\mathbf{w} = [w_1, w_2, \dots, w_D]$



Each of these weights have a simple interpretation:  $w_d$  is the “weight” or contribution of the  $d^{th}$  feature in making this prediction

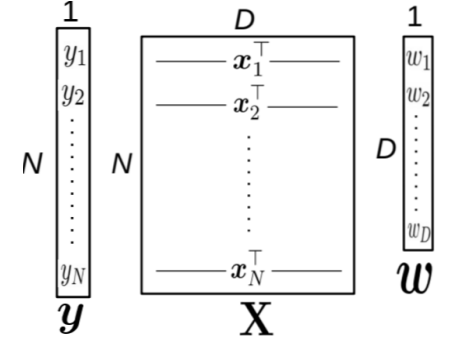
- This simple model can be used for **Linear Regression**
- This simple model can also be used as a “building block” for more complex models, e.g.,
  - Classification (binary/multiclass/multi-output/multi-label) and various other ML/deep learning models
  - Unsupervised learning problems (e.g., dimensionality reduction models)

The “optimal” weights are unknown and have to be learned by solving an **optimization problem**, using some **training data**



# 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  $(\mathbf{x}_n, y_n)$



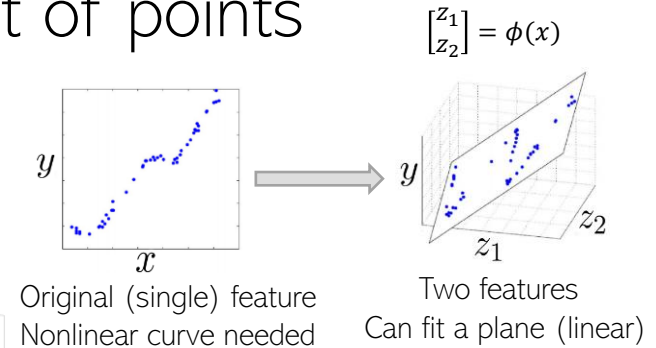
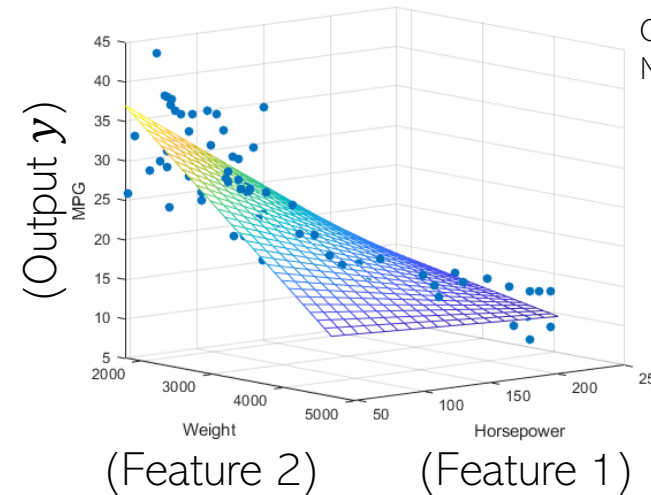
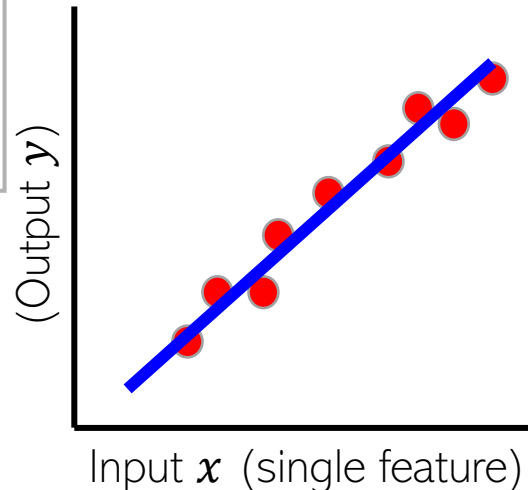


# Linear Regression: Pictorially

- 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

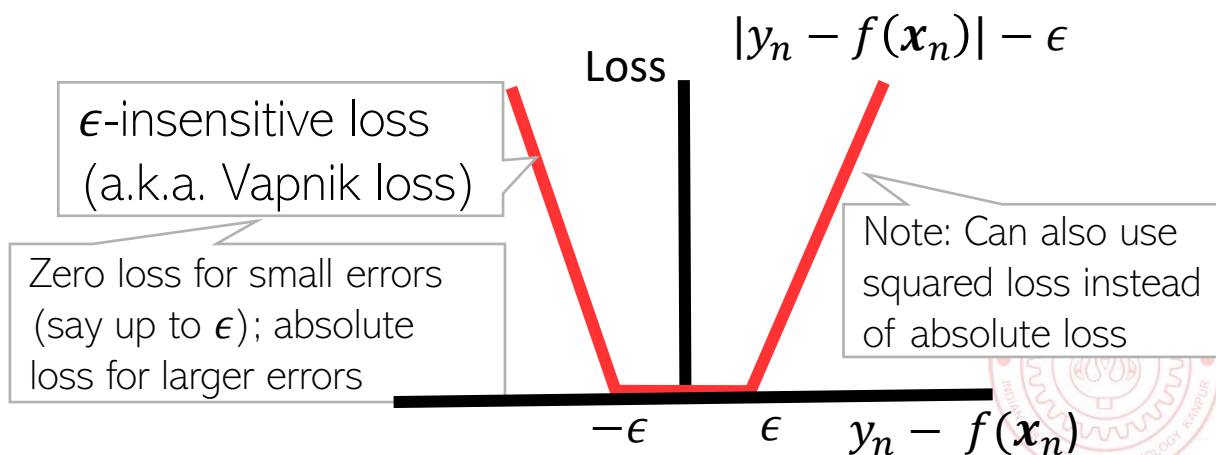
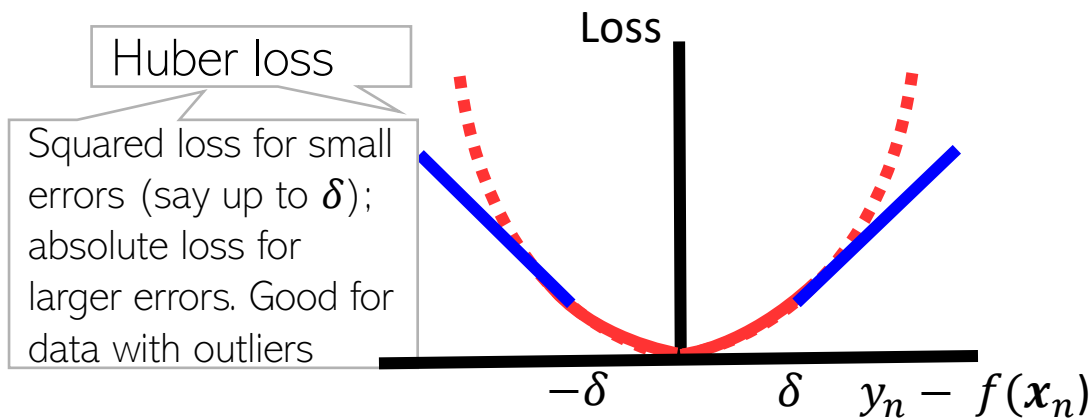
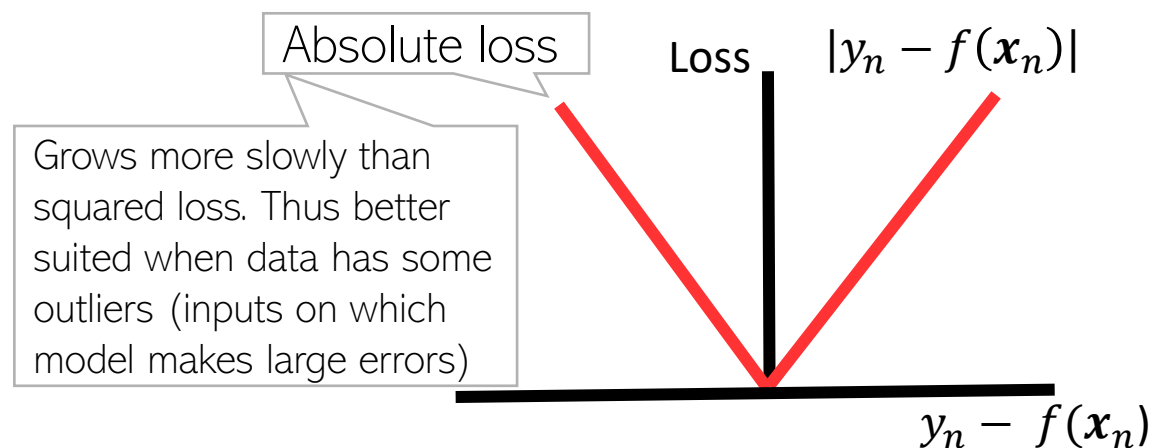
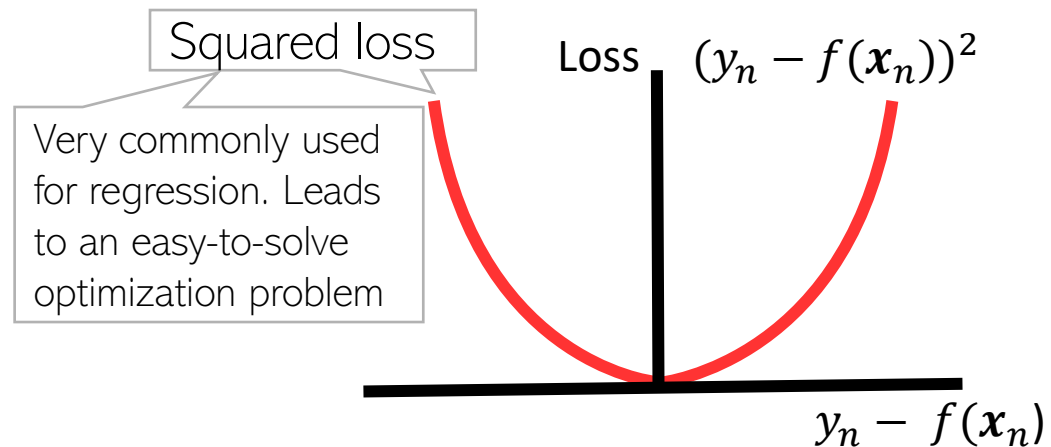
- The line/plane must also predict outputs of the unseen (test) inputs well

# Loss Functions for Regression

- Many possible loss functions for regression problems

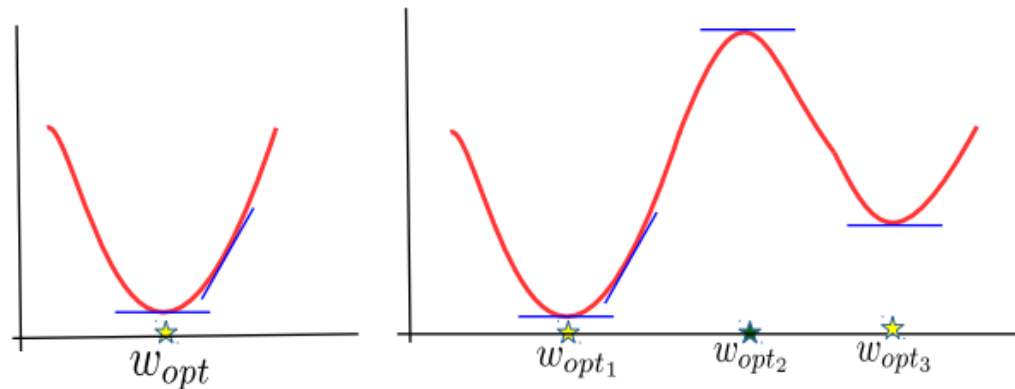
Choice of loss function usually depends on the nature of the data. Also, some loss functions result in easier optimization problem than others

10



# Minimizing Loss Func using First-Order Optimality<sup>11</sup>

- Use basic calculus to find the minima



Called “first order” since only gradient is used and gradient provides the first order info about the function being optimized



The approach works only for very simple problems where the objective is convex and there are no constraints on the values  $\mathbf{w}$  can take

- First order optimality: The gradient  $\mathbf{g}$  must be equal to zero at the optima

$$\mathbf{g} = \nabla_{\mathbf{w}}[L(\mathbf{w})] = \mathbf{0}$$

- Sometimes, setting  $\mathbf{g} = \mathbf{0}$  and solving for  $\mathbf{w}$  gives a closed form solution
- If closed form solution is not available, the gradient vector  $\mathbf{g}$  can still be used in iterative optimization algos, like [gradient descent \(GD\)](#)
  - Note: Even if closed-form solution is possible, GD can sometimes be more efficient



# Minimizing Loss Func. using Iterative Optimization<sup>12</sup>



Can I used this approach to solve **maximization** problems?

For max. problems we can use gradient **ascent**  
 $\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} + \eta_t \mathbf{g}^{(t)}$

Iterative since it requires several steps/iterations to find the optimal solution



Good initialization needed for non-convex functions

For convex functions, GD will converge to the global minima

Will move in the direction of the gradient

**Fact:** Gradient gives the direction of **steepest change** in function's value

## Gradient Descent

- Initialize  $\mathbf{w}$  as  $\mathbf{w}^{(0)}$
- For iteration  $t = 0, 1, 2, \dots$  (or until convergence)
  - Calculate the gradient  $\mathbf{g}^{(t)}$  using the current iterates  $\mathbf{w}^{(t)}$
  - Set the learning rate  $\eta_t$
  - Move in the **opposite** direction of gradient

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta_t \mathbf{g}^{(t)}$$

The learning rate very imp. Should be set carefully (fixed or chosen adaptively). Will discuss some strategies later

Sometimes may be tricky to to assess convergence. Will see some methods later

# 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
- Let's use first order optimality

The “least squares” (LS) problem  
Gauss-Legendre, 18<sup>th</sup> 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)<sup>14</sup>

- 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)<sup>16</sup>

- 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 ( $\ell_2$  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

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

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

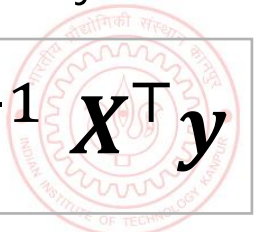
Look at the form of the solution. We are adding a small value  $\lambda$  to the diagonals of the DxD matrix  $\mathbf{X}^\top \mathbf{X}$  (like adding a ridge/mountain to some land)

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

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



Why is the method called "ridge" regression





# Other Ways to Control Overfitting

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

$\ell_1$  norm regularizer

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

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

$\ell_0$  norm regularizer (counts number of nonzeros in  $\mathbf{w}$ )



When should I use these regularizers instead of the  $\ell_2$  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](#)



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)

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

# Gradient Descent for Linear/Ridge Regression

- Just use the GD algorithm with the gradient expressions we derived ☺
- Iterative updates for linear regression will be of the form

$$\begin{aligned}\mathbf{w}^{(t+1)} &= \mathbf{w}^{(t)} - \eta_t \mathbf{g}^{(t)} \\ &= \mathbf{w}^{(t)} + \eta_t \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n) \mathbf{x}_n\end{aligned}$$

Unlike the closed form solution of least squares regression, here we have iterative updates but do not require the expensive matrix inversion of the  $D \times D$  matrix  $\mathbf{X}^\top \mathbf{X}$

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

