

# Linear Models and Learning via Optimization

Piyush Rai

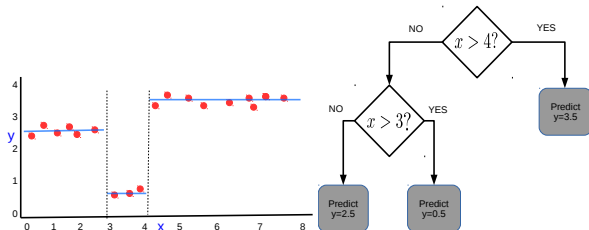
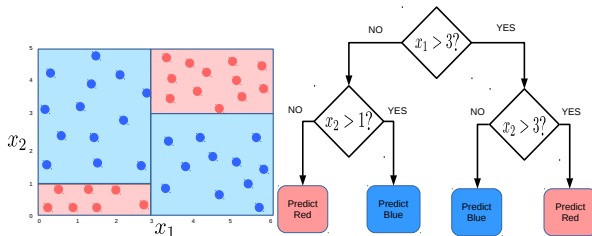
Introduction to Machine Learning (CS771A)

August 9, 2018



# Recap

Decision Trees: Learning by asking questions. Ask the “important” questions first!

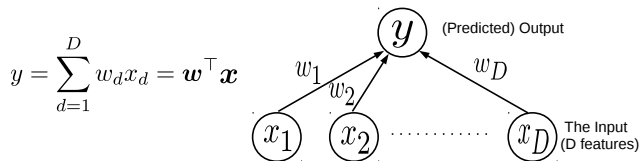


# Linear Models



# Linear Models

- Consider learning to map an input  $\mathbf{x} \in \mathbb{R}^D$  to its output  $y$  (say real-valued)
- Assume the output to be a **linear weighted combination** of the  $D$  input features

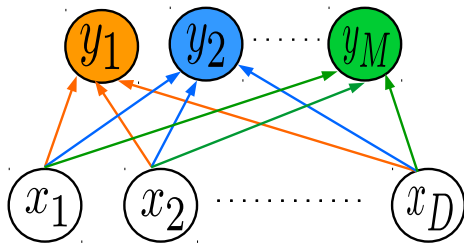


- This is an example of a **linear model** with  $D$  parameters  $\mathbf{w} = [w_1, w_2, \dots, w_D]$
- Inspired by **linear models of neurons**
- $\mathbf{w} \in \mathbb{R}^D$  is also known as the **weight vector**
- Here  $w_d$  denotes how important the  $d$ -th input feature is for predicting  $y$
- The above is basically a linear model for simple **regression** (single, real-valued output  $y$ )
- This basic model can also be used as **building blocks** in many more complex models



# Linear Models for Multi-output Regression

- Can assume each of the  $M$  outputs in  $\mathbf{y} \in \mathbb{R}^M$  to be modeled by a linear model



- Each output  $y_m$  ( $m = 1, \dots, M$ ) modeled by a weight vector  $\mathbf{w}_m \in \mathbb{R}^D$ :  $y_m = \mathbf{w}_m^\top \mathbf{x}$
- The entire model for all  $M$  outputs can be represented as  $\mathbf{y} = \mathbf{W}^\top \mathbf{x}$
- $\mathbf{W} = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_M]$  is a  $D \times M$  matrix



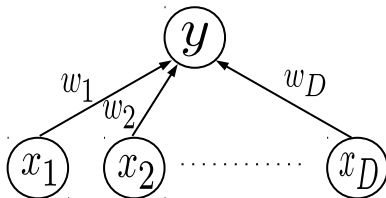
# Linear Models for Binary Classification

- Use the sign of the “score”  $\mathbf{w}^\top \mathbf{x}$  to do predict binary label

$$y = +1 \quad \text{if} \quad \mathbf{w}^\top \mathbf{x} \geq 0$$

$$y = -1 \quad \text{if} \quad \mathbf{w}^\top \mathbf{x} < 0$$

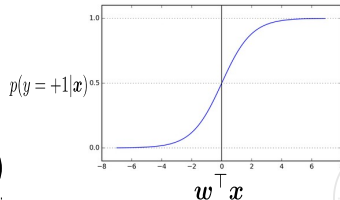
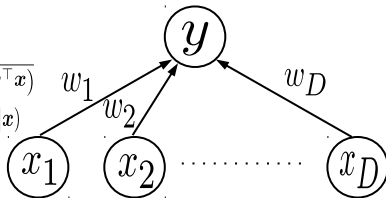
$$y = \text{sign}(\mathbf{w}^\top \mathbf{x})$$



- If desired, can turn the score  $\mathbf{w}^\top \mathbf{x}$  into the probability of the label being +1 (**logistic regression**)

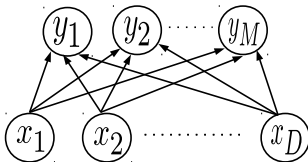
$$p(y = +1|\mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})}$$

$$p(y = -1|\mathbf{x}) = 1 - p(y = +1|\mathbf{x})$$



# Linear Models for Multi-class/Multi-label Classification

- Recall that, in multi-class/multi-label classification,  $\mathbf{y} = [y_1, y_2, \dots, y_M]$  is a vector of length  $M$



- Just like multi-output regression, each component  $y_m$  of  $\mathbf{y}$  can be modeled by a weight vector  $\mathbf{w}_m$
- Need a way to convert  $\mathbf{y} \in \mathbb{R}^M$  to one-hot (for multi-class)/binary vector (for multi-label)
- Note: In some cases, the score need not be converted, e.g.,
  - Can use the index of largest entry in  $\mathbf{y}$  as the predicted class in multi-class classification

0.25	0.6	0.1	0.4	0.2
------	-----	-----	-----	-----

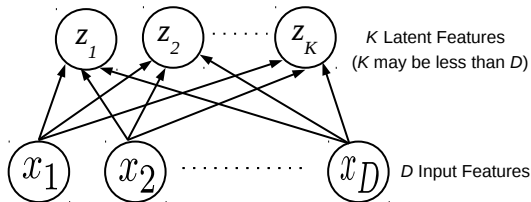
- Can use the indices of top few entries in  $\mathbf{y}$  as the predicted labels in multi-label classification

0.25	0.6	0.1	0.4	0.2
------	-----	-----	-----	-----



# Linear Models for Dimensionality Reduction

- Linear models can be used to reduce data-dimensionality (e.g., [Principal Component Analysis](#))



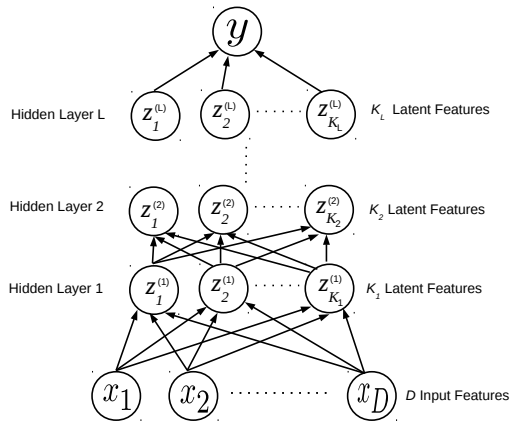
- Note that it looks similar to multi-output regression but the output vector  $\mathbf{z}$  is latent
  - An example of an [unsupervised](#) learning problem
- Need to learn both  $\mathbf{z}$  and  $\mathbf{W}$  in these problems





# Linear Models to construct Deep Neural Networks

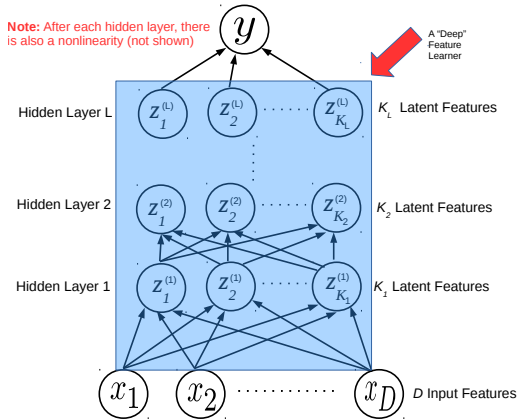
- Linear models are used as basic components of deep neural networks (nonlinear models)



- Each hidden layer has a **learned latent features based representation** of the original input  $\mathbf{x}$ .

# Linear Models to construct Deep Neural Networks

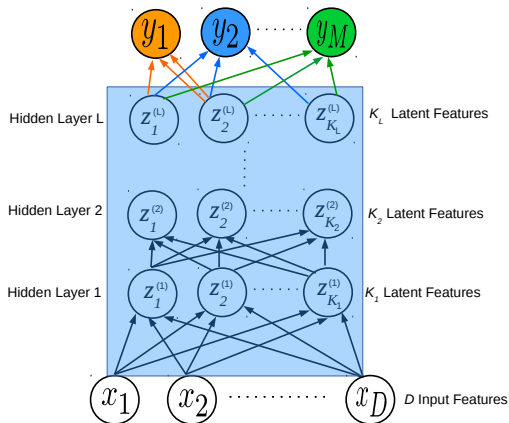
- Linear models are used as basic components of deep neural networks (nonlinear models)



- Each hidden layer has a **learned latent features based representation** of the original input  $\mathbf{x}$

# Linear Models to construct Deep Neural Networks

- Can even construct multiple-output versions of deep neural network

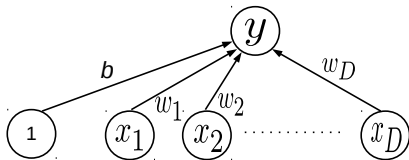


- These can be used for multi-output regression, multi-class/multi-label classification, etc.



# Linear Models with Offset (Bias) Parameter

- Some linear models use an additional bias parameter  $b$

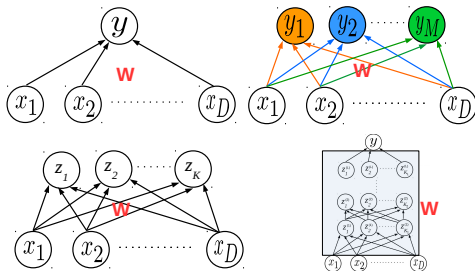


$$y = \sum_{d=1}^D w_d x_d + b = \mathbf{w}^\top \mathbf{x} + b$$

- Can append a constant feature “1” for each input and rewrite as  $y = \mathbf{w}^\top \mathbf{x}$ , with  $\mathbf{x}, \mathbf{w} \in \mathbb{R}^{D+1}$
- We will assume the same and omit the explicit bias for simplicity of notation



# Learning Linear Models



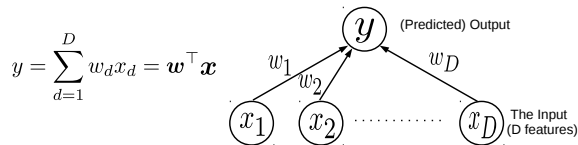
Linear Models are ubiquitous!  
How do we learn them from data?

For linear models, learning = Learning the model parameters (the weights)

We will formulate learning as an **optimization problem** w.r.t. these parameters

# Learning a Linear Model for Regression

- Let's focus on learning the simplest linear model for now: **Linear Regression**



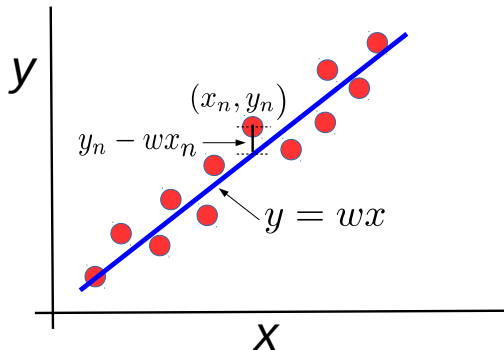
- Suppose we are given regression training data  $\{(\mathbf{x}_n, y_n)\}_{n=1}^N$  with  $\mathbf{x}_n \in \mathbb{R}^D$ , and  $y_n \in \mathbb{R}$
- Let's model the training data using  $\mathbf{w}$  and assume  $y_n \approx \mathbf{w}^\top \mathbf{x}_n, \forall n$  (equivalently  $\mathbf{y} \approx \mathbf{X}\mathbf{w}$ )

The diagram shows the matrix equation  $\mathbf{y} \approx \mathbf{X}\mathbf{w}$ . The vector  $\mathbf{y}$  is on the left, with elements  $y_1, y_2, \dots, y_N$  and a vertical dimension of  $N$ . The matrix  $\mathbf{X}$  is in the middle, with rows  $x_1^\top, x_2^\top, \dots, x_N^\top$  and a horizontal dimension of  $D$ . The vector  $\mathbf{w}$  is on the right, with elements  $w_1, w_2, \dots, w_D$  and a vertical dimension of  $D$ . The approximation symbol  $\approx$  is between  $\mathbf{y}$  and  $\mathbf{X}$ .



# Linear Regression: Pictorially

- With one-dimensional inputs, linear regression would look like



- Error of the model for an example =  $y_n - \mathbf{w}^\top \mathbf{x}_n$  ( $= y_n - wx_n$  for scalar input case)



# Linear Regression

- Define the total error or “loss” on the training data, when using  $\mathbf{w}$  as our model, as

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

- Note: Squared loss chosen for simplicity. Can define other type of losses too (more on this later)
- The best  $\mathbf{w}$  will be the one that minimizes the above error (requires [optimization](#) w.r.t.  $\mathbf{w}$ )

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

- This is known as “[least squares](#)” linear regression (Gauss/Legendre, early 18th century)
- Taking derivative (gradient) of  $\mathcal{L}(\mathbf{w})$  w.r.t.  $\mathbf{w}$  and setting to zero

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

- Simplifying further, we get a [closed form solution](#) for  $\mathbf{w} \in \mathbb{R}^D$

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





# Linear Regression

$$\begin{matrix} 1 \\ y_1 \\ y_2 \\ \vdots \\ y_N \\ \mathbf{y} \end{matrix} \approx \begin{matrix} D \\ \text{---} x_1^{\top} \text{---} \\ \text{---} x_2^{\top} \text{---} \\ \vdots \\ \text{---} x_N^{\top} \text{---} \\ \mathbf{X} \end{matrix} \begin{matrix} 1 \\ w_1 \\ w_2 \\ \vdots \\ w_D \\ \mathbf{w} \end{matrix}$$

- Consider the closed form solution we obtained for linear regression based on least squares

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

- The above closed form solution is nice but has some issues
  - The  $D \times D$  matrix  $\mathbf{X}^{\top} \mathbf{X}$  may not be invertible
  - Based **solely on minimizing the training error**  $\sum_{n=1}^N (y_n - \mathbf{w}^{\top} \mathbf{x}_n)^2 \Rightarrow$  can **overfit** the training data
  - Expensive inversion** for large  $D$ : Can use **iterative optimization** techniques (will come to this later)



# Regularized Linear Regression (a.k.a. Ridge Regression)

- Consider **regularized loss**: Training error +  $\ell_2$ -squared norm of  $\mathbf{w}$ , i.e.,  $\|\mathbf{w}\|_2^2 = \mathbf{w}^\top \mathbf{w} = \sum_{d=1}^D w_d^2$

$$\mathcal{L}_{reg}(\mathbf{w}) = \left[ \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2 + \frac{\lambda}{2} \mathbf{w}^\top \mathbf{w} \right]$$

- Minimizing the above objective w.r.t.  $\mathbf{w}$  does two things
  - Keeps the training error small
  - Keeps the  $\ell_2$  norm of  $\mathbf{w}$  small (and thus also the individual components of  $\mathbf{w}$ ): **Regularization**
- There is a trade-off between the two terms: The **regularization hyperparam**  $\lambda > 0$  controls it
  - Very small  $\lambda$  means almost no regularization (can overfit)
  - Very large  $\lambda$  means very high regularization (can underfit - high training error)
  - Can use cross-validation to choose the “right”  $\lambda$
- The solution to the above optimization problem is:  $\mathbf{w} = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_D)^{-1} \mathbf{X}^\top \mathbf{y}$
- Note that, in this case, regularization also made inversion possible (note the  $\lambda \mathbf{I}_D$  term)



# How $\ell_2$ Regularization Helps Here?

- We saw that  $\ell_2$  regularization encourages the individual weights in  $\mathbf{w}$  to be small
- Small weights ensure that the function  $y = f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$  is **smooth** (i.e., we expect similar  $\mathbf{x}$ 's to have similar  $y$ 's). Below is an informal justification:
- Consider two points  $\mathbf{x}_n \in \mathbb{R}^D$  and  $\mathbf{x}_m \in \mathbb{R}^D$  that are exactly similar in all features **except the  $d$ -th feature** where they differ by a small value, say  $\epsilon$
- Assuming a simple/smooth function  $f(\mathbf{x})$ ,  $y_n$  and  $y_m$  should also be close
- However, as per the model  $y = f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$ ,  $y_n$  and  $y_m$  will differ by  $\epsilon w_d$
- Unless we constrain  $w_d$  to have a small value, the difference  $\epsilon w_d$  would also be very large (which isn't what we want).
- That's why regularizing (via  $\ell_2$  regularization) and making the individual components of the weight vector small helps



# Regularization: Some Comments

- Many ways to regularize ML models (for linear as well as other models)
- Some are based on adding a norm of  $\mathbf{w}$  to the loss function (as we already saw)
  - Using  $\ell_2$  norm in the loss function promotes the individual entries  $\mathbf{w}$  to be small (we saw that)
  - Using  $\ell_0$  norm encourages very few non-zero entries in  $\mathbf{w}$  (thereby promoting “sparse”  $\mathbf{w}$ )

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

- Optimizing with  $\ell_0$  is difficult (NP-hard problem); can use  $\ell_1$  norm as an approximation

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

- **Note:** Since they learn a sparse  $\mathbf{w}$ ,  $\ell_0$  or  $\ell_1$  regularization is also useful for doing **feature selection** ( $w_d = 0$  means feature  $d$  is irrelevant). We will revisit  $\ell_1$  later to formally see why  $\ell_1$  gives sparsity
- Other techniques for regularization: **Early stopping** (of training), **“dropout”**, etc (popular in deep neural networks; will revisit these later when discussing deep learning)



# Linear/Ridge Regression via Gradient Descent

- Both least squares regression and ridge regression require **matrix inversion**

$$\text{Least Squares } \mathbf{w} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}, \quad \text{Ridge } \mathbf{w} = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_D)^{-1} \mathbf{X}^\top \mathbf{y}$$

- Can be **computationally expensive** when  $D$  is very large
- A faster way is to use **iterative optimization**, such as batch or stochastic **gradient descent**
- A basic batch gradient-descent based procedure looks like
  - Start with an initial value of  $\mathbf{w} = \mathbf{w}^{(0)}$
  - Update  $\mathbf{w}$  by moving along the gradient of the loss function  $\mathcal{L}$

$$\mathbf{w}^{(t)} = \mathbf{w}^{(t-1)} - \eta \left. \frac{\partial \mathcal{L}}{\partial \mathbf{w}} \right|_{\mathbf{w}=\mathbf{w}^{(t-1)}} \quad \text{where } \eta \text{ is the learning rate}$$

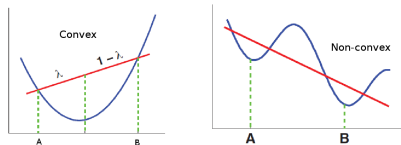
- Repeat until converge
- For least squares, the gradient is  $\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = -\sum_{n=1}^N \mathbf{x}_n (y_n - \mathbf{x}_n^\top \mathbf{w})$  (no matrix inversion involved)
- Such iterative methods for optimizing loss functions are widely used in ML. **Will revisit these later**



# Linear Regression via Gradient-based Methods: Some Notes

We will revisit gradient based methods later but a few things to keep in mind

- Gradient Descent guaranteed to converge to a local minima
- Gradient Descent converges to global minima if the function is **convex**



- A function is convex if second derivative is non-negative everywhere (for scalar functions) or if Hessian is positive semi-definite (for vector-valued functions). For a convex function, every local minima is also a global minima.
- Note: The **squared loss function** in linear regression is **convex**
  - With  $\ell_2$  regularizer, it becomes strictly convex (single global minima)
- For Gradient Descent, the learning rate is important (should not be too large or too small)



# Linear Regression as Solving System of Linear Equations

- Solving  $\mathbf{y} = \mathbf{X}\mathbf{w}$  for  $\mathbf{w}$  is like solving for  $D$  unknowns  $w_1, \dots, w_D$  using  $N$  equations

$$y_1 = x_{11}w_1 + x_{12}w_2 + \dots + x_{1D}w_D$$

$$y_2 = x_{21}w_1 + x_{22}w_2 + \dots + x_{2D}w_D$$

$$\vdots$$

$$y_N = x_{N1}w_1 + x_{N2}w_2 + \dots + x_{ND}w_D$$

- Can therefore view the linear regression problem as a **system of linear equations**
- However, in linear regression, we would rarely have  $N = D$ , but  $N > D$  or  $D > N$
- $N > D$  case is an **overdetermined** system of linear equations ( $\#$  equations  $>$   $\#$  unknowns)
- $D > N$  case is an **underdetermined** system of linear equations ( $\#$  unknowns  $>$   $\#$  equations)
- Thus methods to solve over/underdetermined systems can be used to solve linear regression as well
  - Many of these don't require a matrix inversion (will provide a separate note with details)



# Linear Regression: Some Other Comments

- A simple and interpretable method. Very widely used.
- Least squares and ridge regression are one of the very few ML problems with closed form solutions

$$\text{Least Squares } \mathbf{w} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}, \quad \text{Ridge } \mathbf{w} = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_D)^{-1} \mathbf{X}^\top \mathbf{y}$$

- Many ML problems can be easily reduced to the form  $\mathbf{y} = \mathbf{X}\mathbf{w}$  or  $\mathbf{Y} = \mathbf{X}\mathbf{W}$
- Equivalence to **over/underdetermined** system of linear equations enables us to use efficient solvers (a lot of work in the numerical linear algebra community to scale up linear systems solvers)
  - **An interesting bit:** Note that  $\mathbf{w} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \Rightarrow \mathbf{A}\mathbf{w} = \mathbf{b}$  where  $\mathbf{A} = \mathbf{X}^\top \mathbf{X}$  and  $\mathbf{b} = \mathbf{X}^\top \mathbf{y}$
  - Using the above relation, can solve for  $\mathbf{w}$  by solving  $\mathbf{A}\mathbf{w} = \mathbf{b}$ . A **standard** linear system with  $D$  equations and  $D$  unknowns; can be solved using efficient linear systems solvers.
- The basic (regularized) linear regression can also be easily extended to
  - **Nonlinear Regression**  $y_n \approx \mathbf{w}^\top \phi(\mathbf{x}_n)$  by replacing the original feature vector  $\mathbf{x}_n$  by a nonlinear transformation  $\phi(\mathbf{x}_n)$  (where  $\phi$  may be pre-defined or itself learned)
  - **Generalized Linear Model**  $y_n = g(\mathbf{w}^\top \mathbf{x}_n)$  when response  $y_n$  is **not real-valued** but binary/categorical/count, etc, and  $g$  is a “link function”





# General Supervised Learning as Optimization

- We saw that regularized least squares regression required solving

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \mathcal{L}_{reg}(\mathbf{w}) = \arg \min_{\mathbf{w}} \left[ \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2 + \frac{\lambda}{2} \mathbf{w}^\top \mathbf{w} \right]$$

- This is essentially the **training loss** (called “**empirical loss**”), plus the **regularization** term
- In general, for supervised learning, the goal is to learn a function  $f$ , s.t.  $f(\mathbf{x}_n) \approx y_n, \forall n$
- Moreover, we also want to have a simple  $f$ , i.e., have some regularization
- Therefore, learning the best  $f$  amounts to solving the following **optimization problem**

$$\hat{f} = \arg \min_f \mathcal{L}_{reg}(f) = \arg \min_f \sum_{n=1}^N \ell(y_n, f(\mathbf{x}_n)) + \lambda R(f)$$

where  $\ell(y_n, f(\mathbf{x}_n))$  measures the model  $f$ 's **training loss** on  $(\mathbf{x}_n, y_n)$  and  $R(f)$  is a **regularizer**

- For least squares regression,  $f(\mathbf{x}_n) = \mathbf{w}^\top \mathbf{x}_n$ , and  $R(f) = \mathbf{w}^\top \mathbf{w}$ , and  $\ell(y_n, f(\mathbf{x}_n)) = (y_n - \mathbf{w}^\top \mathbf{x}_n)^2$
- As we'll see later, different supervised learning problems differ in the **choice of  $f$ ,  $R(\cdot)$ , and  $\ell$**

# General Unsupervised Learning as Optimization

- Can we formulate unsupervised learning problems as optimization problems? Yes, of course! :-)
- Consider an unsupervised learning problem with  $N$  inputs  $\mathbf{X} = \{\mathbf{x}_n\}_{n=1}^N$
- Unsupervised, so no labels. Suppose we are interested in learning a new representation  $\mathbf{Z} = \{\mathbf{z}_n\}_{n=1}^N$
- Assume a function  $f$  that models the relationship between  $\mathbf{x}_n$  and  $\mathbf{z}_n$

$$\mathbf{x}_n \approx f(\mathbf{z}_n) \quad \forall n$$

- In this case, we can define a loss function  $\ell(\mathbf{x}_n, f(\mathbf{z}_n))$  that measures how well  $f$  can “reconstruct” the original  $\mathbf{x}_n$  from its new representation  $\mathbf{z}_n$
- This generic unsup. learning problem can thus be written as the following optimization problem

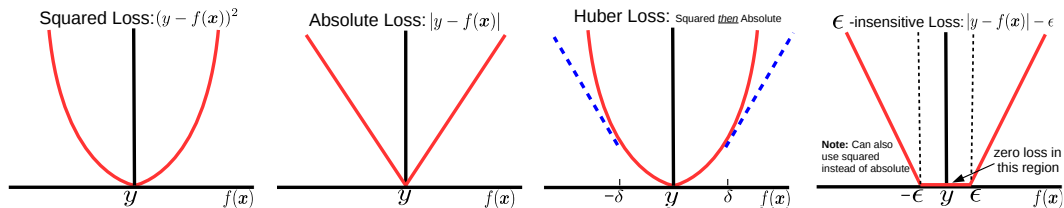
$$\hat{f} = \arg \min_{f, \mathbf{Z}} \sum_{n=1}^N \ell(\mathbf{x}_n, f(\mathbf{z}_n)) + \lambda R(f, \mathbf{Z})$$

- In this case both  $f$  and  $\mathbf{Z}$  need to be learned. Typically learned via alternating optimization



# A Brief Detour: Some Loss Functions

- Some popular loss functions for regression problems<sup>1</sup>



- Absolute/Huber loss preferred if there are outliers in the data
  - Less affected by large errors  $|y - f(x)|$  as compared to the squared loss
- Overall objective function = loss func + some regularizer (e.g.,  $\ell_2$ ,  $\ell_1$ ), as we saw for ridge reg.
- Some objectives easy to optimize (convex and differentiable), some not so (e.g., non-differentiable)
- Will revisit many of these aspects when we talk about optimization techniques for ML

<sup>1</sup> will look at loss functions for classification later when discussing classification in detail

