

COMP 3105 Introduction to Machine Learning

Lecture 1: Linear Regression

Junfeng Wen

Carleton University

Fall 2025

Regression Examples: Estimate Real Numbers

Estimate fish weight

- ▶ Shape (length/width) \mapsto weight

Annual production of a corn farm

- ▶ Rainfall/sunshine/pest levels \mapsto production

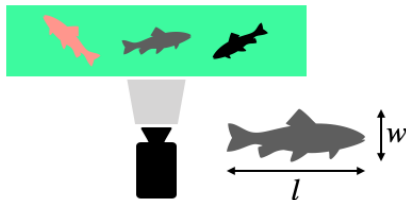
Glucose level of patients with diabetes

- ▶ Insulin injection, patient info \mapsto glucose level

Market price of a new product

- ▶ Cost/quality/transportation \mapsto price

An Example: Estimate Fish Weight



Every fish is measured by [length, width]

- ▶ E.g., [70 cm, 18 cm]
- ▶ Captured by camera on the conveyor belt

Want to estimate its *weight*

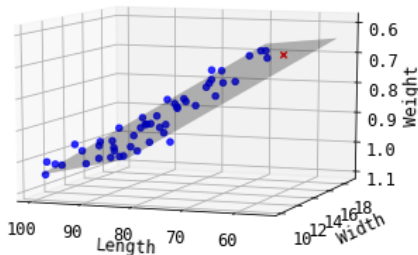
An Example: Estimate Fish Weight (Cont.)

One possible formulation

- ▶ Represent each fish as a (**column**) vector $\mathbf{x} = \begin{bmatrix} 1 \\ 70 \\ 18 \end{bmatrix}$ & $y = 0.7$
- ▶ Plot what we know (blue points; three axes)
- ▶ Find a “plane” (with **augmented** \mathbf{x}):

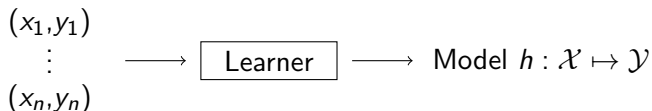
$$h(\mathbf{x}) = w_0 + w_1x_1 + w_2x_2 = \underset{1 \times 3}{\mathbf{x}}^T \underset{3 \times 1}{\mathbf{w}}$$

- ▶ For new fish \mathbf{x}_0 , predict $\hat{y}_0 = h(\mathbf{x}_0)$



General Framework: Supervised Learning

Training data



x_i input in domain \mathcal{X} (e.g., length and width \mathbb{R}^2)

- ▶ A.k.a. features, attributes, variables...
- ▶ Numeric (Length: 0.2, 1.3...)
- ▶ Categorical (Blood type: A/B/AB/O)
- ▶ Ordinal (Difficulty: easy/normal/hard; Discrete but rankable)

y_i output in range \mathcal{Y} (e.g., weight \mathbb{R})

- ▶ A.k.a. ground-truth label (or simply label), target variable...
- ▶ Supervised: training data have both x and y

Goal: Find h that predicts well for novel test $x \in \mathcal{X}$

The I.I.D. Assumption

ML requires certain assumption to work

Assume $\{(x_i, y_i)\}_{i=1}^n$ are independently and identically distributed

- ▶ Independently: $\forall i, (x_i, y_i)$ is freshly (independently) drawn from a probability distribution $P(x, y)$
- ▶ Identical: All examples are drawn from the same $P(x, y)$

The same applies to the novel test data (drawn from $P(x, y)$)

Linear Regression Formulation

Suppose

- ▶ Input vector $\mathbf{x}_i \in \mathcal{X} = \mathbb{R}^d$
- ▶ Output scalar $y_i \in \mathcal{Y} = \mathbb{R}$ (regression)
- ▶ Linear model $h_{\mathbf{w}}(\mathbf{x}) = \mathbf{x}^\top \mathbf{w}$ for some *parameters* $\mathbf{w} \in \mathbb{R}^d$

Goal

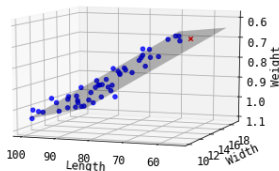
- ▶ Find a good \mathbf{w} so that $\hat{y} \triangleq h_{\mathbf{w}}(\mathbf{x})$ is “close” to y

Closeness? Loss function $L(\hat{y}, y)$

- ▶ Absolute loss (L_1 loss) $L_1(\hat{y}, y) = |\hat{y} - y|$
- ▶ Squared loss (L_2 loss) $L_2(\hat{y}, y) = \frac{1}{2}(\hat{y} - y)^2$

Objective function: Empirical risk minimization (ERM)

$$\min_{\mathbf{w} \in \mathbb{R}^d} \quad \frac{1}{n} \sum_{i=1}^n L(\mathbf{x}_i^\top \mathbf{w}, y_i)$$



Matrix/Vector Notations

The input matrix and output vector

$$X \triangleq \begin{bmatrix} - & \mathbf{x}_1^\top & - \\ - & \mathbf{x}_2^\top & - \\ & \vdots & \\ - & \mathbf{x}_n^\top & - \end{bmatrix} \in \mathbb{R}^{n \times d} \text{ and } \mathbf{y} \triangleq \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \in \mathbb{R}^n$$

Then the prediction

$$\hat{\mathbf{y}} = X\mathbf{w} = \begin{bmatrix} - & \mathbf{x}_1^\top & - \\ - & \mathbf{x}_2^\top & - \\ & \vdots & \\ - & \mathbf{x}_n^\top & - \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix} \in \mathbb{R}^n$$

Want $\hat{\mathbf{y}} \approx \mathbf{y}$

Linear Regression Objective

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n L(\mathbf{x}_i^\top \mathbf{w}, y_i)$$

Matrix/Vector notations

$$X \triangleq \begin{bmatrix} - & \mathbf{x}_1^\top & - \\ - & \mathbf{x}_2^\top & - \\ & \vdots & \\ - & \mathbf{x}_n^\top & - \end{bmatrix} \in \mathbb{R}^{n \times d} \text{ and } \mathbf{y} \triangleq \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \in \mathbb{R}^n$$

► Want $\hat{\mathbf{y}} = X\mathbf{w} \approx \mathbf{y}$

With L_2 loss $L_2(\hat{y}_i, y_i) = \frac{1}{2}(\hat{y}_i - y_i)^2 = \frac{1}{2}(\mathbf{x}_i^\top \mathbf{w} - y_i)^2$

► Recall (squared) L_2 norm $\|\mathbf{a}\|_2^2 = a_1^2 + a_2^2 + \cdots + a_n^2 = \mathbf{a}^\top \mathbf{a}$

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{2n} \|X\mathbf{w} - \mathbf{y}\|_2^2$$

How to Solve It? Direct Approach

$$\min_{\mathbf{w} \in \mathbb{R}^d} J(\mathbf{w}) \triangleq \frac{1}{2n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 = \frac{1}{2n} (\mathbf{X}\mathbf{w} - \mathbf{y})^\top (\mathbf{X}\mathbf{w} - \mathbf{y})$$

Note that this objective function J is convex

Compute the gradient

$$\begin{aligned} \nabla J(\mathbf{w}) &= \frac{1}{2n} \nabla \left[(\mathbf{X}\mathbf{w} - \mathbf{y})^\top (\mathbf{X}\mathbf{w} - \mathbf{y}) \right] \quad \text{if } f(u) = (au - b)^2, f'(u)? \text{ (chain rule)} \\ &= \frac{1}{n} (\mathbf{X}\mathbf{w} - \mathbf{y}) \times? \nabla (\mathbf{X}\mathbf{w} - \mathbf{y}) \quad \text{if } f(u) = au - b, f'(u)? \\ &= \frac{1}{n} \left(\underset{n \times d}{\mathbf{X}} \underset{d \times 1}{\mathbf{w}} - \underset{n \times 1}{\mathbf{y}} \right) \times? \underset{n \times d}{\mathbf{X}} \quad \text{Dimension matched?} \\ &= \frac{1}{n} \underset{d \times n}{\mathbf{X}^\top} \underset{n \times 1}{(\mathbf{X}\mathbf{w} - \mathbf{y})} \end{aligned}$$

Grad shares the same dim as the params

Set it to zero to find the best \mathbf{w} . Assume that $\mathbf{X}^\top \mathbf{X}$ is invertible

$$\mathbf{X}^\top (\mathbf{X}\mathbf{w} - \mathbf{y}) = \mathbf{0}_d \iff \mathbf{X}^\top \mathbf{X}\mathbf{w} = \mathbf{X}^\top \mathbf{y} \iff \mathbf{w}^* = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

How to Solve It? Indirect Approach

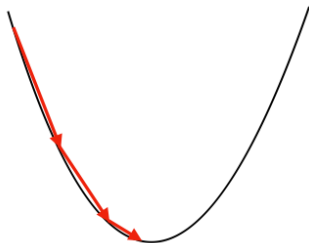
$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{2n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

Direct approach: compute the gradient and set to zero

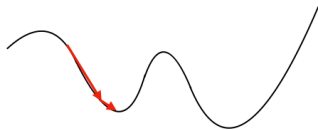
- ▶ Not always possible

Indirect approach: iterative algorithm (gradient descent)

- ▶ Initial $\mathbf{w}^{(0)}$ (e.g. $\mathbf{w}^{(0)} = \mathbf{0}_d$, a vector of all zeros)
- ▶ $\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta \nabla J(\mathbf{w}^{(t)})$ with step size $\eta > 0$
- ▶ That's why convexity could be important



Convex

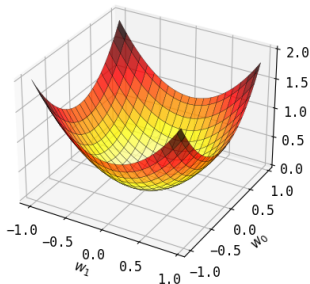


Non-convex

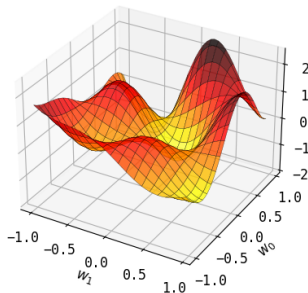
How to Solve It? Indirect Approach

Gradient descent 3D visualization

- $\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta \nabla J(\mathbf{w}^{(t)})$ with step size $\eta > 0$



Convex

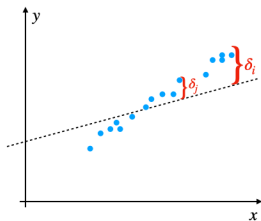


Non-convex

How about Other Losses? L_1 Loss

L_1 loss minimization

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n |\mathbf{x}_i^\top \mathbf{w} - y_i|$$



- ▶ Convex but not *smooth* (abs not differentiable at 0)
- ▶ Direct approach is not possible

Other solvers?

- ▶ For point (\mathbf{x}_i, y_i) , defined a “tolerance” $\delta_i \geq 0$

$$\begin{aligned} |\hat{y}_i - y_i| \leq \delta_i &\iff -\delta_i \leq \hat{y}_i - y_i \leq \delta_i \\ &\iff \mathbf{x}_i^\top \mathbf{w} - y_i \leq \delta_i \text{ and } y_i - \mathbf{x}_i^\top \mathbf{w} \leq \delta_i \end{aligned}$$

- ▶ Minimize sum of tolerances $\min_{\mathbf{w}, \delta_i} \frac{1}{n} \sum_{i=1}^n \delta_i$

L_1 Loss Minimization in Matrix Form

Using matrix/vector notations

- ▶ Recall L_1 norm $\|\mathbf{a}\|_1 = |a_1| + |a_2| + \cdots + |a_n|$

$$\min_{\mathbf{w} \in \mathbb{R}^d} \quad \frac{1}{n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_1$$

Let $\boldsymbol{\delta} \in \mathbb{R}_+^n$ be the “tolerance” between $\hat{\mathbf{y}}$ and \mathbf{y}

- ▶ Always non-negative $\boldsymbol{\delta} \succeq \mathbf{0}_n$

$$\min_{\mathbf{w}, \boldsymbol{\delta}} \quad \boldsymbol{\delta}^\top \mathbf{1}_n = \sum_{i=1}^n \delta_i \quad \text{Want tolerance sum to be small}$$

$$\text{s.t.} \quad \boldsymbol{\delta} \succeq \mathbf{0}_n$$

$$\mathbf{X}\mathbf{w} - \mathbf{y} \preceq \boldsymbol{\delta} \quad \text{Sandwich the difference}$$

$$\mathbf{y} - \mathbf{X}\mathbf{w} \preceq \boldsymbol{\delta} \quad \text{on both sides} \Leftrightarrow |\mathbf{X}\mathbf{w} - \mathbf{y}| \preceq \boldsymbol{\delta}$$

Linear programming (LP) problem (Why? Linear in both \mathbf{w} and $\boldsymbol{\delta}$), which can be solved efficiently

How about Other Losses? L_∞ Loss

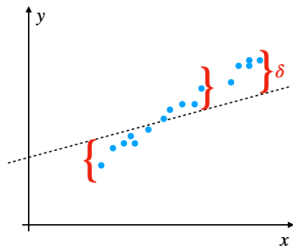
In general, the L_p loss: $\|X\mathbf{w} - \mathbf{y}\|_p^p, 0 \leq p \leq \infty$

L_∞ loss: $\|X\mathbf{w} - \mathbf{y}\|_\infty \triangleq \max_{i=1,\dots,n} |\mathbf{x}_i^\top \mathbf{w} - y_i|$

► Control the maximum gap (worse point)

Let $\delta \geq 0$ be the scalar “tolerance” (such that $|\mathbf{x}_i^\top \mathbf{w} - y_i| \leq \delta$)

$$\begin{aligned} \min_{\mathbf{w}, \delta} \quad & \delta \\ \text{s.t.} \quad & \delta \geq 0 \\ & X\mathbf{w} - \mathbf{y} \preceq \delta \cdot \mathbf{1}_n \\ & \mathbf{y} - X\mathbf{w} \preceq \delta \cdot \mathbf{1}_n \end{aligned}$$

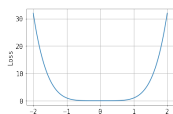


Again, an LP problem

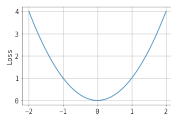
Which Loss to Use?

Different p values give different models \mathbf{w}^*

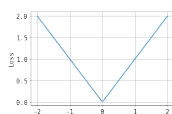
- ▶ $1 \leq p \leq \infty$: objective function is convex
- ▶ $1 < p < \infty$: objective function is smooth and differentiable
- ▶ $0 < p < 1$: *not* convex, but grow slowly as \hat{y} deviates from y
 - ▶ More robust, not affected much by outliers



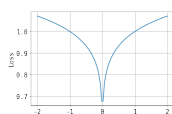
$p = 5$



$p = 2$



$p = 1$

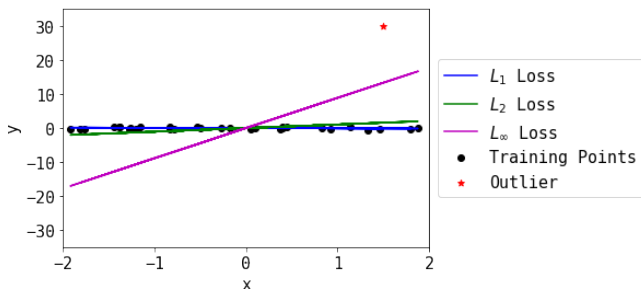


$p = 0.1$

Loss with different p values for one data point (x-axis: $\hat{y} - y$)

Which Loss to Use?

An example: Learning with outliers



- ▶ Large p (L_∞) chases the outlier (focus on the worse error)
- ▶ Small p (L_1) is robust (treat every error equally)

Always use small p ? No, because...

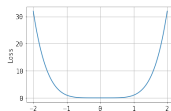
- ▶ $p < 1$ is not convex and difficult to minimize

Which Loss to Use?

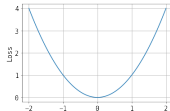
No free lunch

Pros and cons of convex versus robust losses

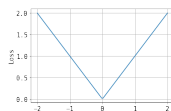
Loss type	Pros 👍	Cons 👎
Convex loss $p \in [1, \infty]$	Efficient optimization	Sensitive to outlier
Robust loss $p \in [0, 1)$	Slow growth, robust	Difficult to optimize



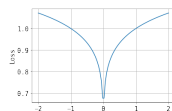
$p = 5$



$p = 2$



$p = 1$



$p = 0.1$

Recap

This lecture: linear regression models $\mathbf{y} \approx X\mathbf{w}$

- ▶ L_2 loss \rightarrow analytic solution
- ▶ L_1 & L_∞ losses \rightarrow linear programming (LP)
- ▶ Different p : convex versus robust losses

References

Bishop and Nasrabadi (2006, Sec.3.1)

Hastie et al. (2009, Sec.2.3.1 & Sec.3.2)

Bishop, C. M. and Nasrabadi, N. M. (2006). *Pattern recognition and machine learning*, volume 4. Springer.

Hastie, T., Tibshirani, R., Friedman, J. H., and Friedman, J. H. (2009). *The elements of statistical learning: data mining, inference, and prediction*, volume 2. Springer.

Appendix: CVXOPT LP Formulation for L_1 Loss

$$\begin{array}{ll}\min_{\mathbf{w}, \delta} & \delta^\top \mathbf{1}_n \\ \text{s.t.} & \delta \succeq \mathbf{0}_n \\ & X\mathbf{w} - \mathbf{y} \preceq \delta \\ & \mathbf{y} - X\mathbf{w} \preceq \delta\end{array}\qquad \begin{array}{ll}\min_{\mathbf{x}} & \mathbf{c}^\top \mathbf{x} \\ \text{s.t.} & G\mathbf{x} \preceq \mathbf{h}\end{array}$$

\mathbf{x} in CVXOPT means unknowns (not a training input vector)

$$\mathbf{x} = \begin{bmatrix} \mathbf{w} \\ d \times 1 \\ \delta \\ n \times 1 \end{bmatrix} \in \mathbb{R}^{d+n}$$

Want $\mathbf{c}^\top \mathbf{x} = \delta^\top \mathbf{1}_n = \sum_{i=1}^n \delta_i$. What is $\mathbf{c} \in \mathbb{R}^{d+n}$ then?

$$\mathbf{c}^\top \mathbf{x} = \begin{bmatrix} \underbrace{0, 0, \dots, 0}_{d \text{ zeros}}, \underbrace{1, 1, \dots, 1}_{n \text{ ones}} \end{bmatrix} \begin{bmatrix} \mathbf{w} \\ \delta \end{bmatrix} = \delta^\top \mathbf{1}_n$$

Appendix: CVXOPT LP Formulation for L_1 Loss

$$\begin{aligned} \text{s.t. } \delta &\succeq \mathbf{0}_n & \text{s.t. } G\mathbf{x} &\preceq \mathbf{h} \\ X\mathbf{w} - \mathbf{y} &\preceq \delta \\ \mathbf{y} - X\mathbf{w} &\preceq \delta \end{aligned}$$

Constraints in CVXOPT

$$G\mathbf{x} = \begin{bmatrix} - & G_1: & - \\ - & G_2: & - \\ & \vdots & \\ - & G_k: & - \end{bmatrix} \mathbf{x} \preceq \begin{bmatrix} h_1 \\ h_2 \\ \vdots \\ h_k \end{bmatrix} = \mathbf{h}$$

Each row is a constraint

$$G_{i:}\mathbf{x} \leq h_i \quad i = 1, 2, \dots, k$$

Appendix: CVXOPT LP Formulation for L_1 Loss

$$\begin{aligned} \text{s.t. } \quad & \delta \succeq \mathbf{0}_n \iff -\delta \preceq \mathbf{0}_n & \text{s.t. } & G\mathbf{x} \preceq \mathbf{h} \\ & X\mathbf{w} - \mathbf{y} \preceq \delta \\ & \mathbf{y} - X\mathbf{w} \preceq \delta \end{aligned}$$

Three sets of constraints (color coded)

$$\begin{matrix} G \\ 3n \times (d+n) \end{matrix} \cdot \begin{matrix} \mathbf{x} \\ (d+n) \times 1 \end{matrix} = \begin{bmatrix} G^{(1)} \\ G^{(2)} \\ G^{(3)} \\ n \times (d+n) \end{bmatrix} \begin{bmatrix} \mathbf{w} \\ \delta \end{bmatrix} \preceq \begin{bmatrix} \mathbf{h}^{(1)} \\ \mathbf{h}^{(2)} \\ \mathbf{h}^{(3)} \\ n \times 1 \end{bmatrix} = \mathbf{h}$$

Want $G^{(1)} \cdot \begin{bmatrix} \mathbf{w} \\ \delta \end{bmatrix} \preceq \mathbf{h}^{(1)} \iff -\delta \preceq \mathbf{0}_n$. What are $G^{(1)}$ and $\mathbf{h}^{(1)}$?

$$\begin{bmatrix} G^{(11)} & G^{(12)} \\ n \times d & n \times n \end{bmatrix} \cdot \begin{bmatrix} \mathbf{w} \\ \delta \end{bmatrix} = G^{(11)} \mathbf{w} + G^{(12)} \delta = -\delta \preceq \mathbf{0}_n = \mathbf{h}^{(1)}$$

So $G^{(11)} = \mathbf{0}_{n \times d}$, $G^{(12)} = -I_n$ and $\mathbf{h}^{(1)} = \mathbf{0}_n$

Appendix: CVXOPT LP Formulation for L_1 Loss

$$\text{s.t. } \delta \succeq \mathbf{0}_n \iff -\delta \preceq \mathbf{0}_n \quad \text{s.t. } G\mathbf{x} \preceq \mathbf{h}$$

$$X\mathbf{w} - \mathbf{y} \preceq \delta$$

$$\mathbf{y} - X\mathbf{w} \preceq \delta$$

Three sets of constraints (color coded)

$$\begin{matrix} G \\ 3n \times (d+n) \end{matrix} \cdot \begin{matrix} \mathbf{x} \\ (d+n) \times 1 \end{matrix} = \begin{bmatrix} \textcolor{blue}{G}^{(1)} \\ \textcolor{brown}{G}^{(2)} \\ \textcolor{red}{G}^{(3)} \end{bmatrix}_{\substack{n \times (d+n) \\ n \times (d+n) \\ n \times (d+n)}} \begin{bmatrix} \mathbf{w} \\ \delta \end{bmatrix}_{\substack{(d+n) \times 1 \\ n \times 1}} \preceq \begin{bmatrix} \textcolor{blue}{h}^{(1)} \\ \textcolor{brown}{h}^{(2)} \\ \textcolor{red}{h}^{(3)} \end{bmatrix}_{\substack{n \times 1 \\ n \times 1 \\ n \times 1}} = \mathbf{h}$$

$$\text{Want } \textcolor{brown}{G}^{(2)}_{n \times (d+n)} \cdot \begin{bmatrix} \mathbf{w} \\ \delta \end{bmatrix}_{\substack{(d+n) \times 1 \\ n \times 1}} \preceq \textcolor{brown}{h}^{(2)}_{n \times 1} \iff X\mathbf{w} - \delta \preceq \mathbf{y}. \text{ What are } \textcolor{brown}{G}^{(2)}, \textcolor{brown}{h}^{(2)}?$$

$$\begin{bmatrix} \textcolor{brown}{G}^{(21)} & \textcolor{brown}{G}^{(22)} \end{bmatrix}_{\substack{n \times d & n \times n}} \cdot \begin{bmatrix} \mathbf{w} \\ \delta \end{bmatrix}_{\substack{d \times 1 \\ n \times 1}} = \textcolor{brown}{G}^{(21)}_{n \times d} \mathbf{w}_{d \times 1} + \textcolor{brown}{G}^{(22)}_{n \times n} \delta_{n \times 1} = X\mathbf{w} - \delta \preceq \mathbf{y} = \textcolor{brown}{h}^{(2)}$$

$$\text{So } \textcolor{brown}{G}^{(21)} = X, \textcolor{brown}{G}^{(22)} = -I_n \text{ and } \textcolor{brown}{h}^{(2)} = \mathbf{y}$$

Appendix: CVXOPT LP Formulation for L_1 Loss

$$\text{s.t. } \delta \succeq \mathbf{0}_n \iff -\delta \preceq \mathbf{0}_n \quad \text{s.t. } G\mathbf{x} \preceq \mathbf{h}$$

$$X\mathbf{w} - \mathbf{y} \preceq \delta$$

$$\mathbf{y} - X\mathbf{w} \preceq \delta$$

Three sets of constraints (color coded)

$$\begin{matrix} G \\ 3n \times (d+n) \end{matrix} \cdot \begin{matrix} \mathbf{x} \\ (d+n) \times 1 \end{matrix} = \begin{bmatrix} \textcolor{blue}{G}^{(1)} \\ \textcolor{brown}{G}^{(2)} \\ \textcolor{red}{G}^{(3)} \\ n \times (d+n) \end{bmatrix} \begin{bmatrix} \mathbf{w} \\ \delta \end{bmatrix} \preceq \begin{bmatrix} \textcolor{blue}{h}^{(1)} \\ \textcolor{brown}{h}^{(2)} \\ \textcolor{red}{h}^{(3)} \\ n \times 1 \end{bmatrix} = \mathbf{h}$$

$$\text{Want } \begin{matrix} \textcolor{red}{G}^{(3)} \\ n \times (d+n) \end{matrix} \cdot \begin{bmatrix} \mathbf{w} \\ \delta \end{bmatrix} \preceq \begin{matrix} \textcolor{red}{h}^{(3)} \\ n \times 1 \end{matrix} \iff -X\mathbf{w} - \delta \preceq -\mathbf{y}. \text{ Then } \textcolor{red}{G}^{(3)}, \textcolor{red}{h}^{(3)}?$$

$$\begin{bmatrix} \textcolor{red}{G}^{(31)} & \textcolor{red}{G}^{(32)} \\ n \times d & n \times n \end{bmatrix} \cdot \begin{bmatrix} \mathbf{w} \\ \delta \end{bmatrix} = \begin{matrix} \textcolor{red}{G}^{(31)} \\ n \times d \end{matrix} \begin{matrix} \mathbf{w} \\ d \times 1 \end{matrix} + \begin{matrix} \textcolor{red}{G}^{(32)} \\ n \times n \end{matrix} \begin{matrix} \delta \\ n \times 1 \end{matrix} = -X\mathbf{w} - \delta \preceq -\mathbf{y} = \textcolor{red}{h}^{(3)}$$

$$\text{So } \textcolor{red}{G}^{(31)} = -X, \textcolor{red}{G}^{(32)} = -I_n \text{ and } \textcolor{red}{h}^{(3)} = -\mathbf{y}$$

Appendix: CVXOPT LP Formulation for L_1 Loss

To summarize

$$\begin{aligned} \min_{\mathbf{x}} \quad & \underbrace{[0, 0, \dots, 0, 1, 1, \dots, 1]}_{\mathbf{c}^\top} \cdot \underbrace{\begin{bmatrix} \mathbf{w} \\ \delta \end{bmatrix}}_{\mathbf{x}} \\ \text{s.t.} \quad & \underbrace{\begin{bmatrix} \mathbf{0}_{n \times d} & -\mathbf{l}_n \\ \mathbf{X} & -\mathbf{l}_n \\ -\mathbf{X} & -\mathbf{l}_n \end{bmatrix}}_G \cdot \underbrace{\begin{bmatrix} \mathbf{w} \\ \delta \end{bmatrix}}_{\mathbf{x}} \preceq \underbrace{\begin{bmatrix} \mathbf{0}_n \\ \mathbf{y} \\ -\mathbf{y} \end{bmatrix}}_{\mathbf{h}} \end{aligned}$$

After solving $\mathbf{x} = \begin{bmatrix} \mathbf{w} \\ \delta \end{bmatrix}$

- ▶ \mathbf{w} model parameters
- ▶ δ tolerances

Exercise/Assignment: L_∞ loss