# Math 342W Lecture 23

If you see correlation, there will be causation somewhere

Spurious Correlation: Concluding $x$ and $y$ are correlated when they're not
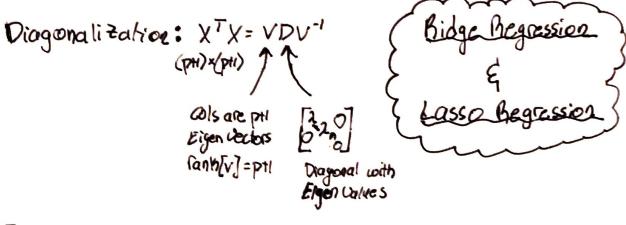
If $\hat{y} = b_0 + b_1 x_1 + \ldots + b_p x_p$, how can we interperet $b_1$

Wrong: Holding the rest of the features constant, a change of 1 in $x_1$ results in $b_1$ change in $y$'s mean.

Correct: When comparing two mutually observed observation (A) and (B) are sampled in the same way as observations in the training set where (A) has a $x_1$ value one unit larger than the $x_1$ of (B), and share same values $x_2, \ldots, x_p$ then (A) is predicted to have a response $y$ that differs by $b_1$ units on average from response of (B), assuming linear model is true.

---

$X = \begin{array}{c} \\ n \end{array}\overset{p}{\boxed{\phantom{XXXX}}}, p > n$

$X^T X$ not invertible now, $\vec{b}_{OLS}$ DNE...

$\Rightarrow X^T X$

$\Rightarrow \vec{b}_{ridge} = (X^T X + \lambda I)^{-1} X^T \vec{y}$

Essentially Cantor's Diagolization plus a little "shift" to add more $n$.

Diagonalization: $X^T X = VDV^{-1}$

$(p+1) \times (p+1)$

Cols are p+1
Eigen vectors
$rank[V] = p+1$

$\begin{bmatrix} \lambda_1 & & 0 \\ & \lambda_2 & \\ 0 & & \lambda_n \end{bmatrix}$

Diagonal with
Eigen values

Ridge Regression
&
Lasso Regression

$\vec{b}_{ridge} = (X^T X + \lambda I)^{-1} X^T \vec{y}$

$(VDV^{-1} + \lambda I)^{-1} = (VDV^{-1} + \lambda I VV^{-1})^{-1} = (VDV^{-1} + V(\lambda I)V^{-1})^{-1}$

$= (V(D + \lambda I)V^{-1})^{-1} = V(D + \lambda I)^{-1}V^{-1} \Rightarrow \begin{bmatrix} \lambda_1 + \lambda & & \\ & \lambda_2 + \lambda & 0 \\ & & \ddots \\ 0 & & \lambda_n + \lambda \\ & & & \lambda \end{bmatrix}$

Now we have shifted $\lambda$'s up to $n$
on the diagonal and $0^+$ $\lambda$'s up to p+1
so now we are full rank.

Consider $\lambda$: $\vec{b}_{ridge} = argmin \{ SSE + \lambda \|\vec{w}\|_2^2 \}$

Consider $\lambda$: $\vec{b}_{lasso} = argmin \{ SSE + \lambda \|\vec{w}\|_1 \}$

$\lambda \in \mathbb{R} > 0$

Regularization
$(\vec{b} \to \vec{0})$

Roughly speaking, Ridge is used for prediction and Lasso
is used for feature selection.

Consider $\lambda$: $\vec{b}_{elastic} = argmin \{ SSE + \lambda (\alpha \|\vec{w}\|_1 + (1-\alpha) \|\vec{w}\|_2^2) \}$