# What does overfitting look like?

- Validation set MSE higher than train set MSE
- Validation set accuracy lower than ~~validation~~ train set accuracy
- regression: predictions not smooth enough —
  fitting noise, not trend

- classification: decision boundary not smooth enough
  — fitting noise, not trend

---

~~Factors~~ Tools we have available to address overfitting:

- # layers, # units per layer
- weight regularization
- Drop out

---

# Process for finding a model:

1) Read the literature or find examples from a similar setting or application

2) Choose your best guess at a good starting point. Fit to training data & evaluate on validation data.

3) Increase model capacity until you are overfitting

4) Regularize model
   - Add L1 or L2 regularization
   - Add drop out
   - Remove layers/units
   - Reduce # of epochs (early stopping)

Can also tune things like learning rate, which optimizer you are using.

5) Refit to combined training & validation data and evaluate on test set.

# L2 regularization in terms of gradient calculations:

$$J(b,w) = \frac{1}{m} \sum_{i=1}^{m} J^{(i)}(b,w) + \cancel{\frac{1}{m}} \sum_{l=1}^{L} \lambda^{[l]} \|w^{[l]}\|_2^2$$

stands for $\sum_i \sum_j (w_{ij}^{[l]})^2$

$$\frac{d}{dw^{[l]}} J(b,w) = \left(\begin{array}{c}\text{stuff from} \\ \text{backpropagation}\end{array}\right) + \lambda^{[l]} \cdot 2 \cdot w^{[l]}$$

$\underbrace{\phantom{xxxxxxxxxx}}$

matrix of same shape as $w^{[l]}$ $(n_{l-1}, n_l)$ where entry $i,j$ is

$$\frac{\partial J(b,w)}{\partial w_{ij}^{[l]}}$$

because $\dfrac{\partial}{\partial w_{ij}^{[l]}} \sum_{l=1}^{L} \lambda^{[l]} \sum_i \sum_j (w_{ij}^{[l]})^2$

$$= 2\lambda^{[l]} w_{ij}^{[l]}$$

---

So a gradient descent update step looks like:

$$w^{[l]} = w^{[l]} - \text{learning-rate} * \left(\begin{array}{c}\text{Stuff from} \\ \text{backpropagation}\end{array} + 2\lambda w^{[l]}\right)$$

$$= w^{[l]}\underbrace{\left(1 - 2\lambda * \text{learning-rate}\right)}_{} - \text{learning-rate} * \left(\begin{array}{c}\text{stuff from} \\ \text{backpropagation}\end{array}\right)$$

every gradient descent
step tries to shrink
weights towards 0
(unless offset by a
gradient term)