
FUNCTION PARAMETERIZATION

Function transformations

Transformation Rules for Functions		
Function Notation	Type of Transformation	Change to Coordinate Point
$f(x) + d$	Vertical translation up d units	$(x, y) \rightarrow (x, y + d)$
$f(x) - d$	Vertical translation down d units	$(x, y) \rightarrow (x, y - d)$
$f(x + c)$	Horizontal translation left c units	$(x, y) \rightarrow (x - c, y)$
$f(x - c)$	Horizontal translation right c units	$(x, y) \rightarrow (x + c, y)$
$-f(x)$	Reflection over x-axis	$(x, y) \rightarrow (x, -y)$
$f(-x)$	Reflection over y-axis	$(x, y) \rightarrow (-x, y)$
$af(x)$	Vertical stretch for $ a > 1$	$(x, y) \rightarrow (x, ay)$
	Vertical compression for $0 < a < 1$	
$f(bx)$	Horizontal compression for $ b > 1$	$(x, y) \rightarrow \left(\frac{x}{b}, y \right)$
	Horizontal stretch for $0 < b < 1$	

<https://www.onlinemathlearning.com/parent-functions.html>

Rigid

Non-Rigid

Parameterizing a function: $y = f(x | p) = Af\left(\frac{x - x_o}{w}\right) + S$

Parameters and Function transformations

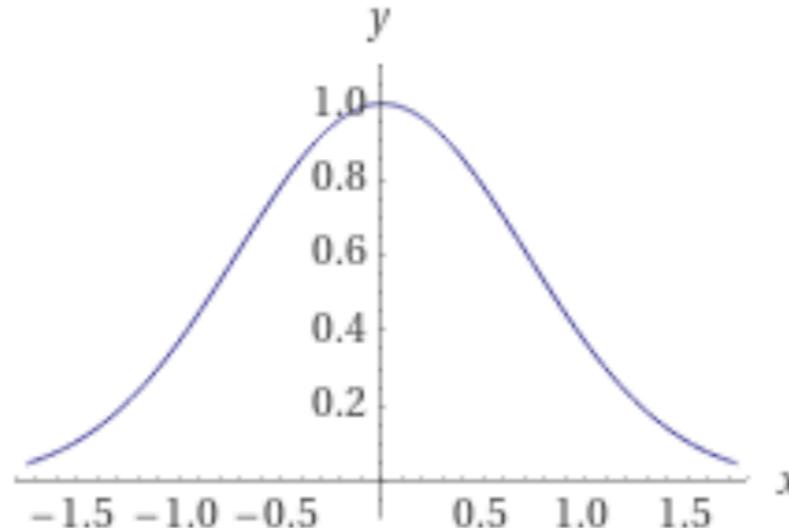
PARAMETRIC MODEL = Functional form + Parameterization

(underlying shape)

(stretching and shifting)

$y = f(x) \rightarrow \text{"parent function"}$

$$f(x) = e^{-x^2}$$



Parameterized form-1:

$$y = f(x | p) = Af\left(\frac{x - x_o}{w}\right) + S$$

$$p = (A, x_o, w, S)$$

A = vertical stretch (amplitude)

x_o = horizontal shift (recentering parameter)

w = horizontal stretch (width parameter)

S = vertical shift (shift parameter)

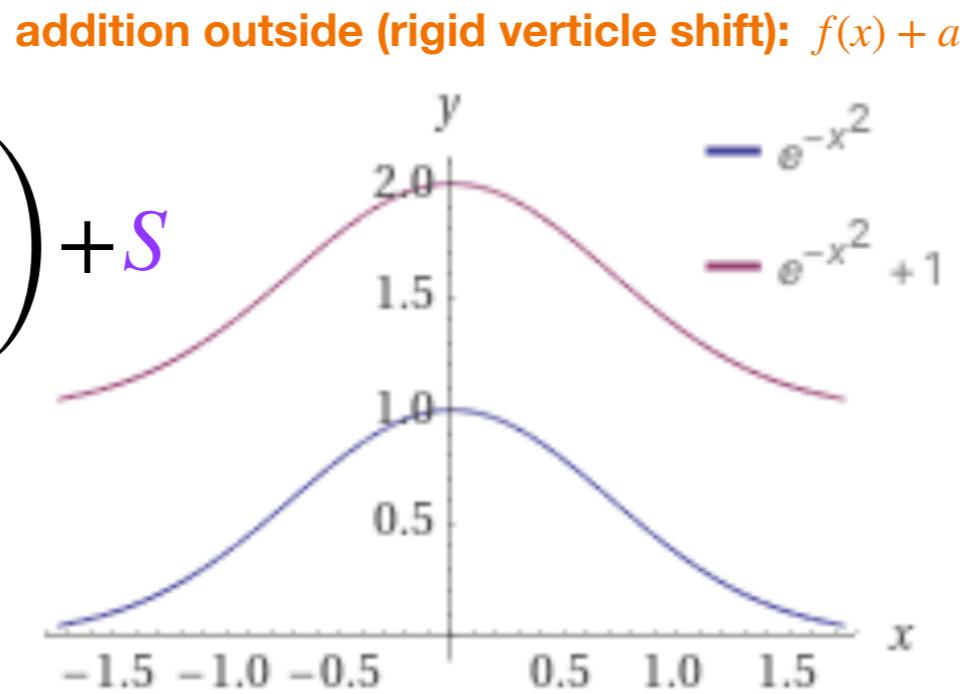
Parameterized form-2: $y = f(x | p) = Af(w_1x + b) + S$

(weights and bias)

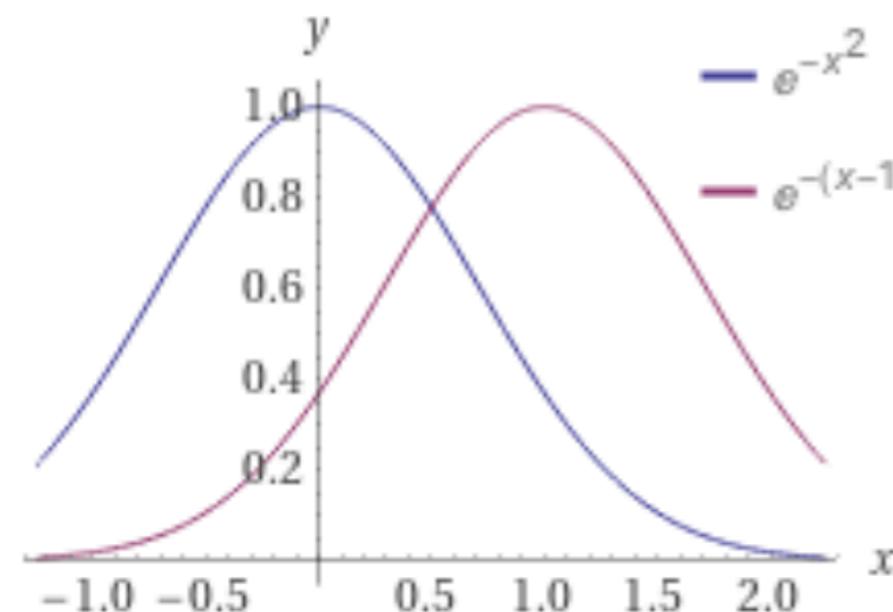
$$w_1 = \frac{1}{w} \quad \text{and} \quad b = -\frac{x_o}{w}$$

Example: Gaussian

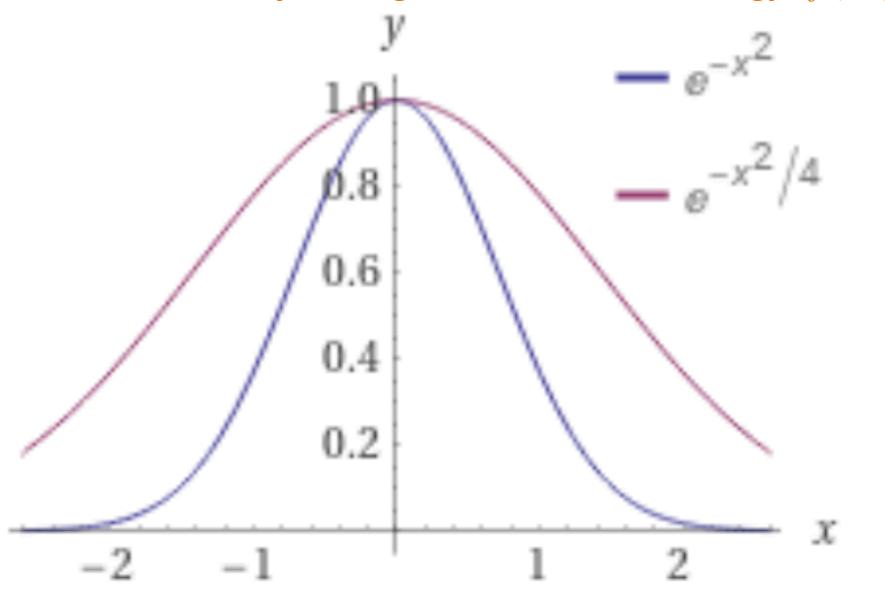
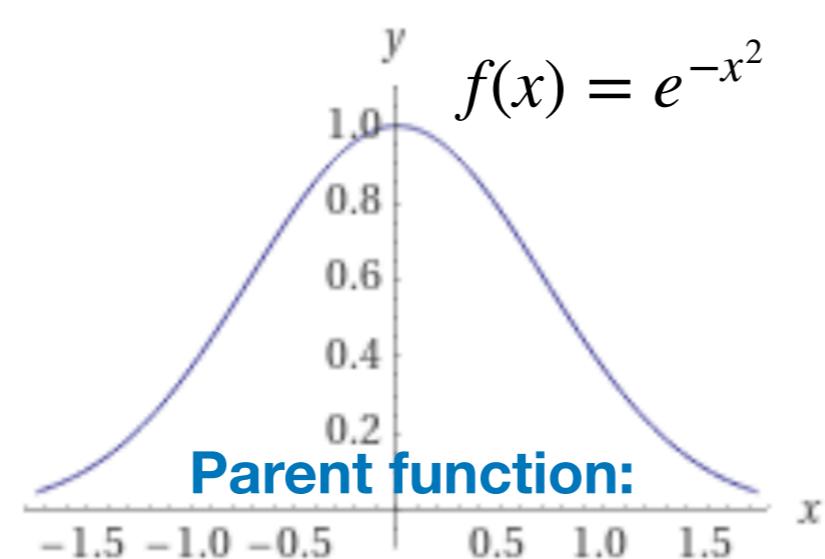
$$y = f(x | \mathbf{p}) = Af\left(\frac{x - x_o}{w}\right) + S$$



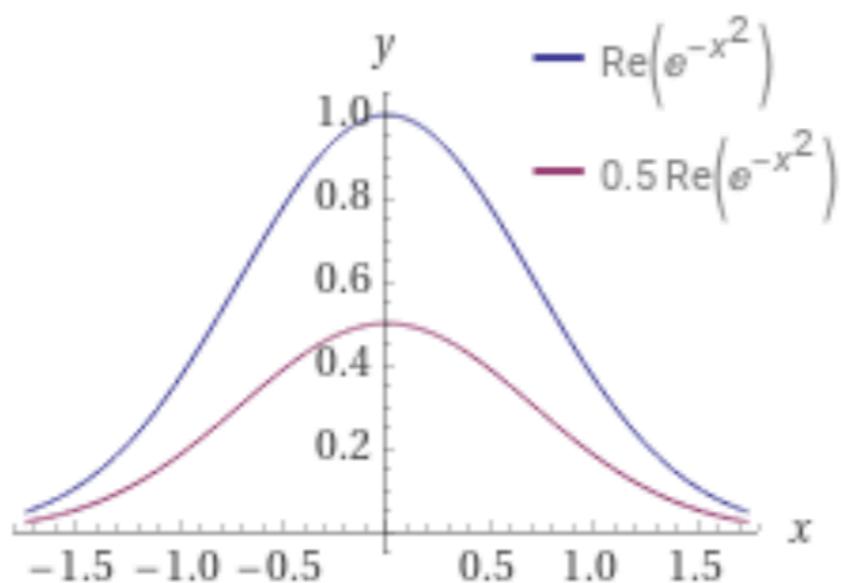
addition inside (rigid horizontal shift): $f(x - a)$



multiplication inside (non-rigid horizontal scaling): $f(ax)$

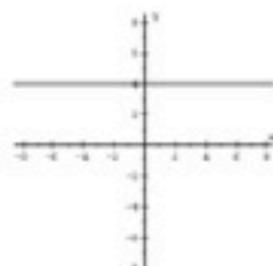
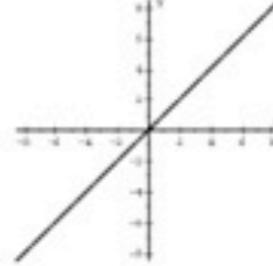
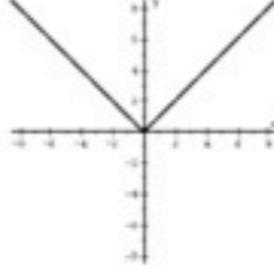
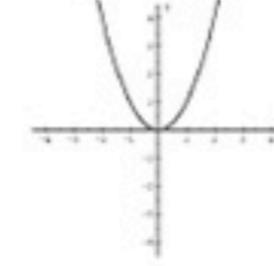
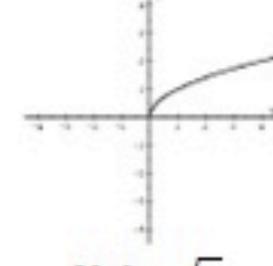
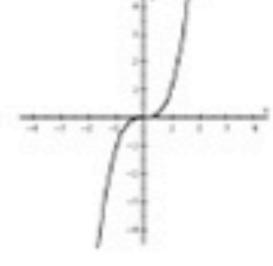
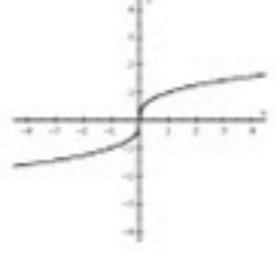
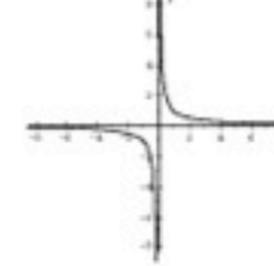
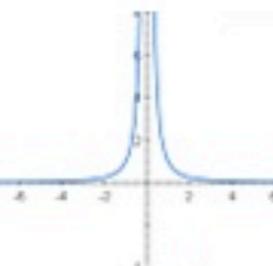
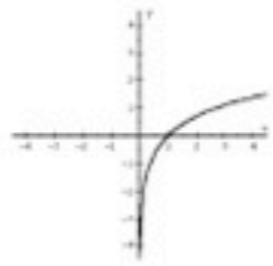
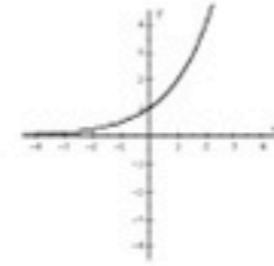
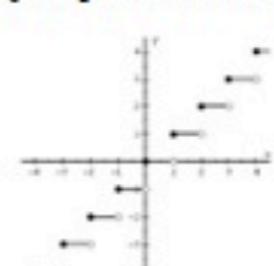
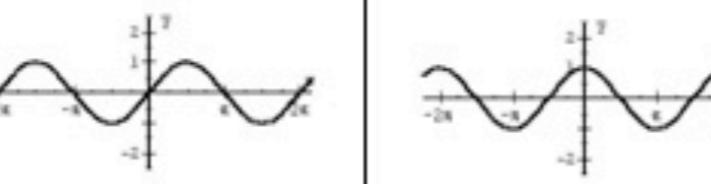
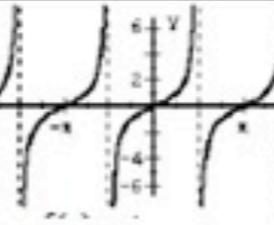


multiplication outside (non-rigid verticle scaling): $af(x)$

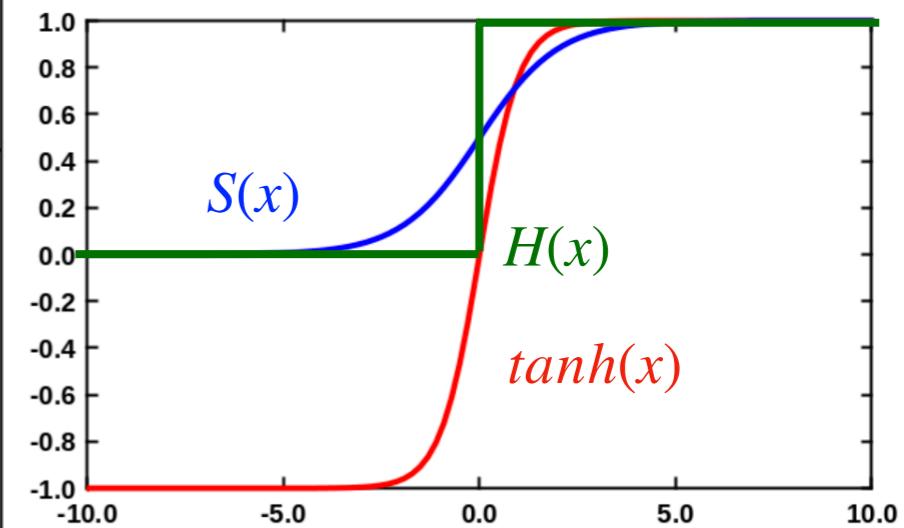


Common parent functions

Model=Functional form + Parameterization

Constant  $f(x) = c$	Linear  $f(x) = x$	Absolute Value  $f(x) = x $	Quadratic  $f(x) = x^2$
Square Root  $f(x) = \sqrt{x}$	Cubic  $f(x) = x^3$	Cube Root  $f(x) = \sqrt[3]{x}$	Reciprocal/Inverse/ Rational  $f(x) = \frac{1}{x}$
Rational  $f(x) = \frac{1}{x^2}$	Logarithmic  $f(x) = \ln(x)$	Exponential  $f(x) = e^x$	Greatest Integer (Step Function)  $f(x) = [[x]]$
Trigonometric Functions →  $f(x) = \sin(x)$ $f(x) = \cos(x)$			 $f(x) = \tan(x)$

IMPORTANT FUNCTIONS IN DEEP LEARNING



Sigmoid (aka. logistic) function

$$S(x) = \frac{1}{1 + e^{-x}} = \frac{e^x}{e^x + 1}$$

Hyperbolic Tangent function

$$\tanh(x) = 2(S(2x) - 0.5)$$

"Heaviside (aka. step) function"

$$H(x) := \begin{cases} 1, & x > 0 \\ 0, & x \leq 0 \end{cases}$$

Basic models

MODEL = Functional form + Parameterization

(underlying shape)

$y = f(x) \rightarrow \text{"parent function"}$

(stretching and shifting)

$\mathbf{p} = (A, x_o, w, S)$

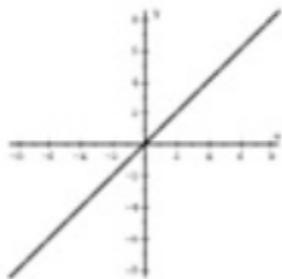
$$y = f(x | \mathbf{p}) = Af\left(\frac{x - x_o}{w}\right) + S$$

Examples:

$$y = f(x | \mathbf{p}) = mx + b$$

Linear

Linear

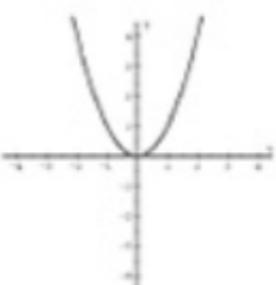


$\mathbf{p} = (m, b)$

$$y = f(x | \mathbf{p}) = Ax^2 + Bx + C$$

Quadratic regressions

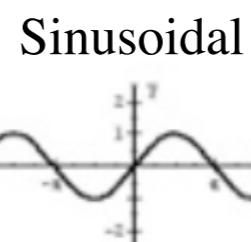
Quadratic



$\mathbf{p} = (A, B, C)$

$$y = f(x | \mathbf{p}) = A_o \sin(\omega t + \phi) + c$$

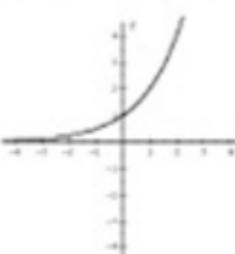
Exponential



Sinusoidal

$\mathbf{p} = (A_o, \omega, \phi, c)$

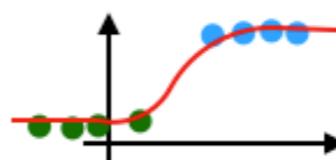
$$y = f(x | \mathbf{p}) = ae^{\frac{(x - x_o)}{w}} + d$$



$\mathbf{p} = (a, w, x_o, d)$

$$y = f(x | \mathbf{p}) = \frac{A}{1 + e^{-\frac{(x - x_o)}{w}}} + S$$

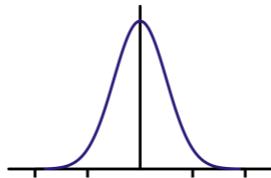
Logistic



$\mathbf{p} = (A, w, x_o, S)$

$$y = f(x | \mathbf{p}) = A \cdot \exp^{-\left(\frac{(x - x_o)}{w}\right)^2} + S$$

Gaussian



$\mathbf{p} = (A, x_o, w, S)$

Higher order models \nearrow **OUTPUTS** $\rightarrow y$ (response)
INPUTS $\rightarrow x$ (variable) \rightarrow **Model=Functional form + Parameterization**

Polynomial regression:

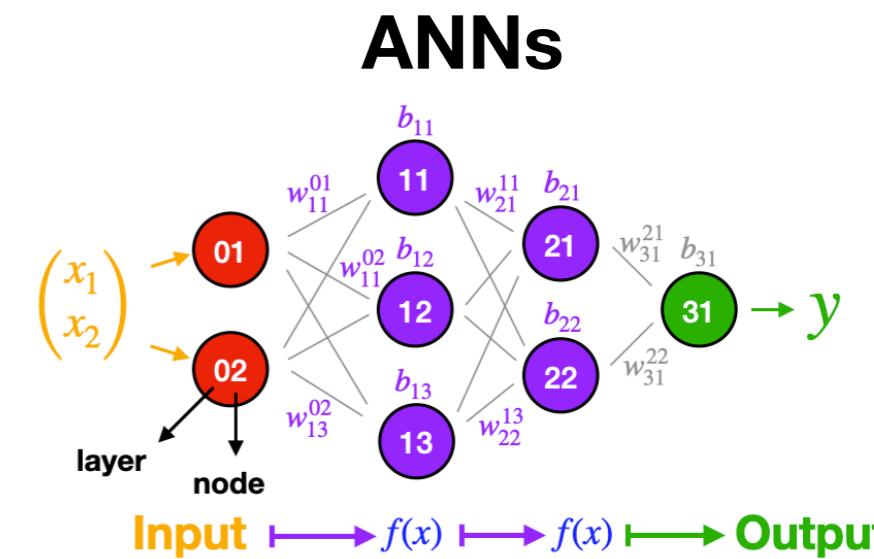
$$y = f(x) = \sum_{n=0}^N A_n x^n$$

$$\mathbf{P} = (A_0, A_1, \dots, A_N)$$

$$y = f(x) = \frac{a_0}{2} + \sum_{n=1}^N \left(a_n \cos\left(\frac{2\pi}{P}nx\right) + b_n \sin\left(\frac{2\pi}{P}nx\right) \right)$$

$$\mathbf{P} = (a_0, b_0, a_1, b_1 \dots a_N,)$$

Fourier series



Neural network model

$$y = ANN(\mathbf{x} | \mathbf{p})$$

fitting parameters

Module-3

More parameters \rightarrow more complex model (more flexible)

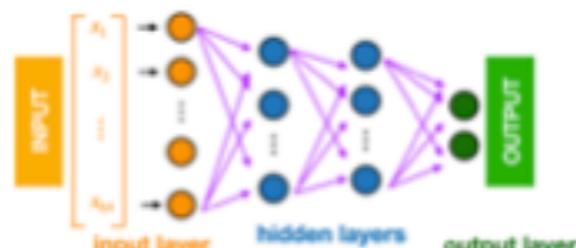
Taylor series

$$f(x) = \sum_{n=0}^{\infty} \frac{f^{(n)}(a)}{n!} (x-a)^n \approx f(a) + \frac{f'(a)}{1!}(x-a) + \frac{f''(a)}{2!}(x-a)^2 + \frac{f'''(a)}{3!}(x-a)^3 + \dots$$

PARAMETRIC MODELING

polynomial regression
logistic regression
mixture models, k-means
hidden Markov models
factor analysis / pPCA / PMF
Linear regression

**Artificial Neural
networks**



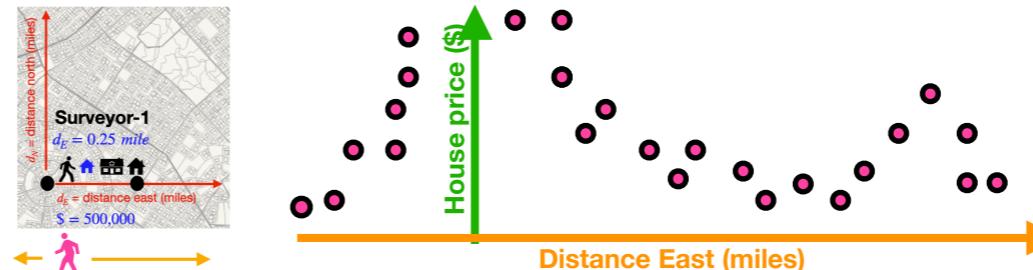
Auto-regression algorithms (timeseries)

General curve fitting

Common ML workflow summary

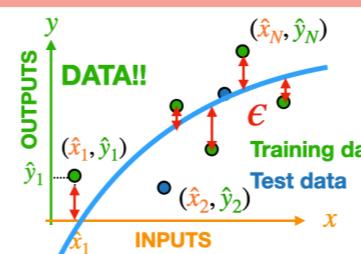
Source and recommended reading: Chollet p111-p116

- Defining the problem and assembling a dataset



- Choosing a measure of success

- Choosing the accuracy/error metrics and a loss function



$$\text{RMSE: } \epsilon(\mathbf{p}) = \left(\frac{\sum_i^N (y(\hat{x}_i | \mathbf{p}) - \hat{y}_i)^2}{N} \right)^{\frac{1}{2}}$$

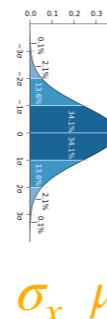
- Deciding on an evaluation protocol



- Preparing your data

Makes it easier to choose parameter initial conditions (IC)

$$\hat{x} = \begin{pmatrix} 0.25 \\ 0.30 \\ \vdots \\ 3.14 \end{pmatrix}$$



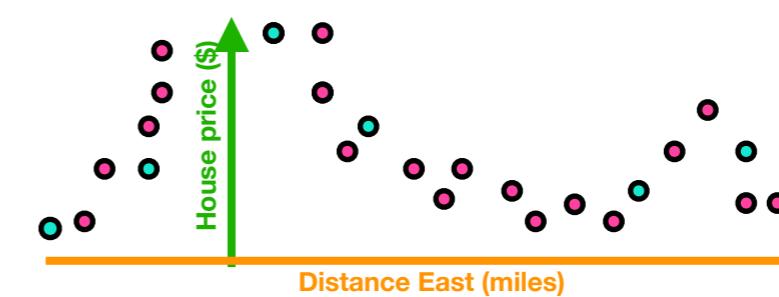
$$x_i \rightarrow \tilde{x}_i = \frac{x_i - \mu_x}{\sigma_x}$$

$$\approx [-1.0, 1.0]$$

$$\hat{y} = \begin{pmatrix} 500,000 \\ 410,000 \\ \vdots \\ 210,000 \end{pmatrix}$$

$$y_i \rightarrow \tilde{y}_i = \frac{y_i - \mu_y}{\sigma_y}$$

$$\approx [-1.0, 1.0]$$



● Train
● Test

- Developing a model that does better than some baseline

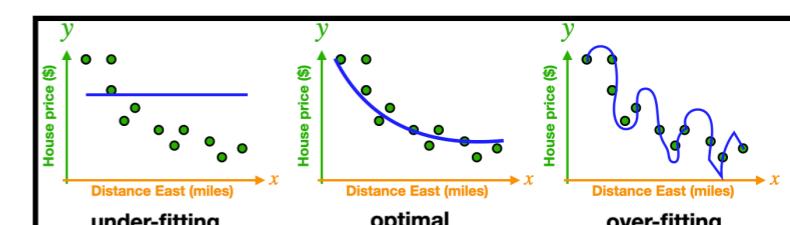
Select a model

$M(x | \mathbf{p}) = \text{Model} = \text{Functional form} + \text{Parameterization}$

$$M(x | \mathbf{p}) = p_1 e^{-\left(\frac{x-p_2}{p_3}\right)^2} + p_4 e^{-\left(\frac{x-p_5}{p_6}\right)^2} + p_7$$

$$\mathbf{p} = (p_0, p_1, p_2, p_3, p_4, p_5, p_6, p_7)$$

- Choosing an optimizer, learning rate, and other hyper-parameters

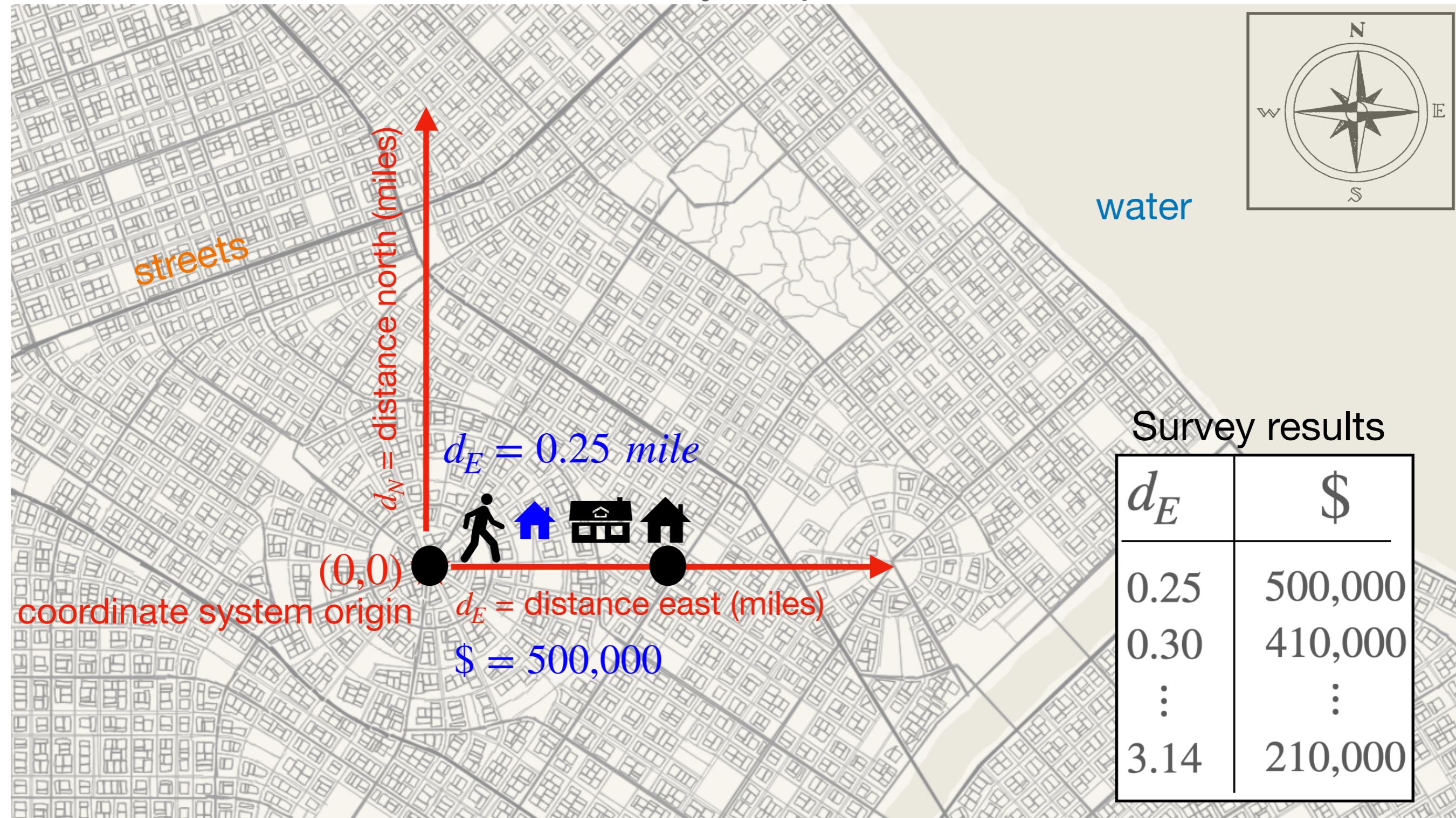


★ Scaling up: developing a model that overfits

★ Regularizing your model and tuning your hyper-parameters

Real Estate Price example

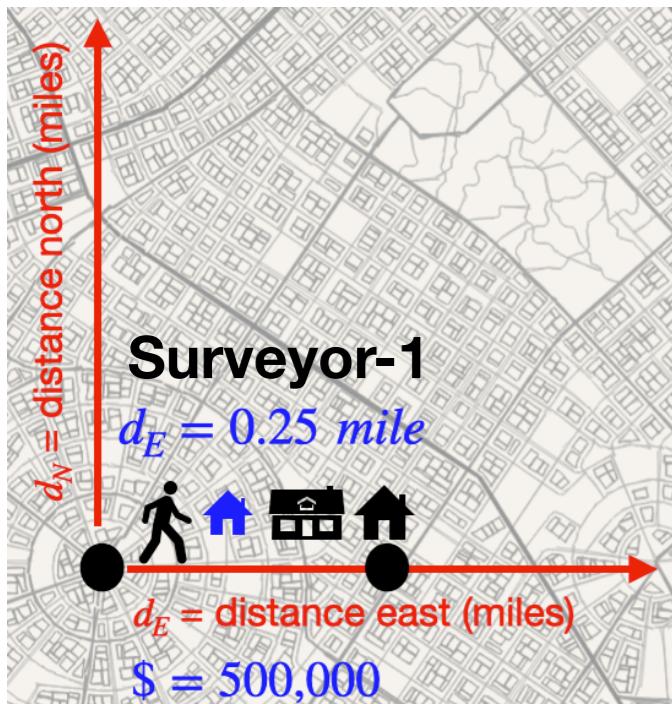
City map



<https://maps.probabilistictrain.com/#/>

$N = 50$

Data



Spreadsheet:

d_E	\$
0.25	500,000
0.30	410,000
:	:
3.14	210,000

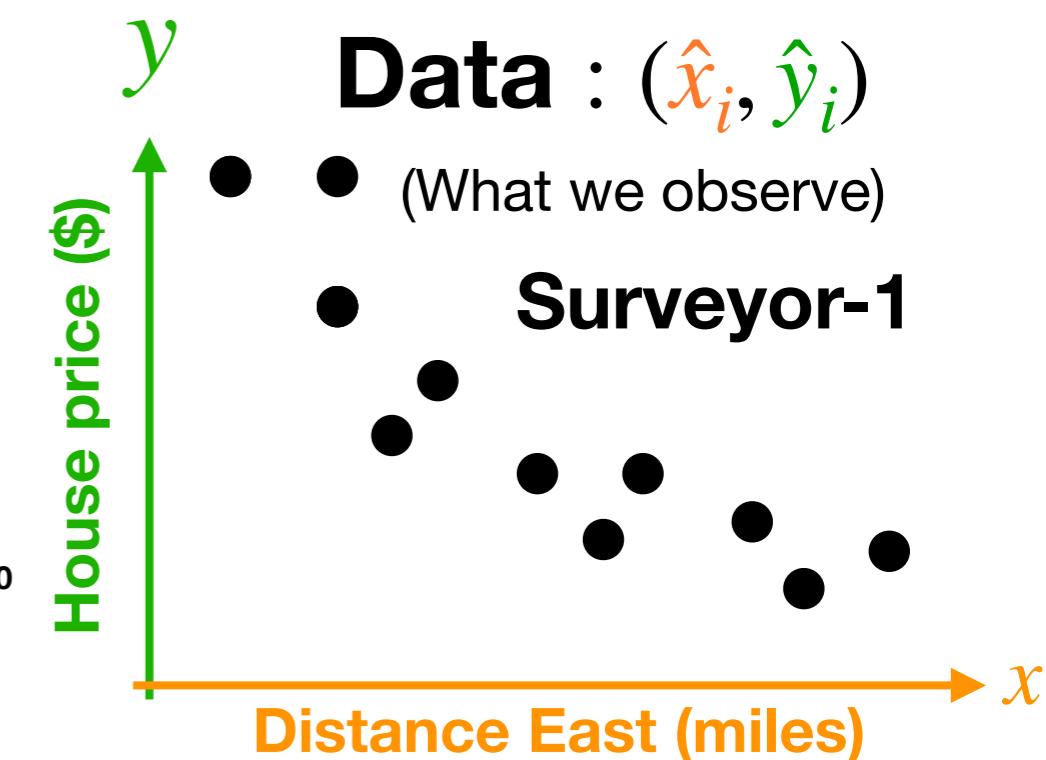
$N = 50$
ordered
Pairs

(\hat{x}_1, \hat{y}_1) Sample-1
 \vdots
 $(\hat{x}_{50}, \hat{y}_{50})$ Sample-50

Inputs: \hat{x}_i Outputs: \hat{y}_i

Hat \rightarrow data

No-Hat \rightarrow continuous variable



"batch=stacked samples"

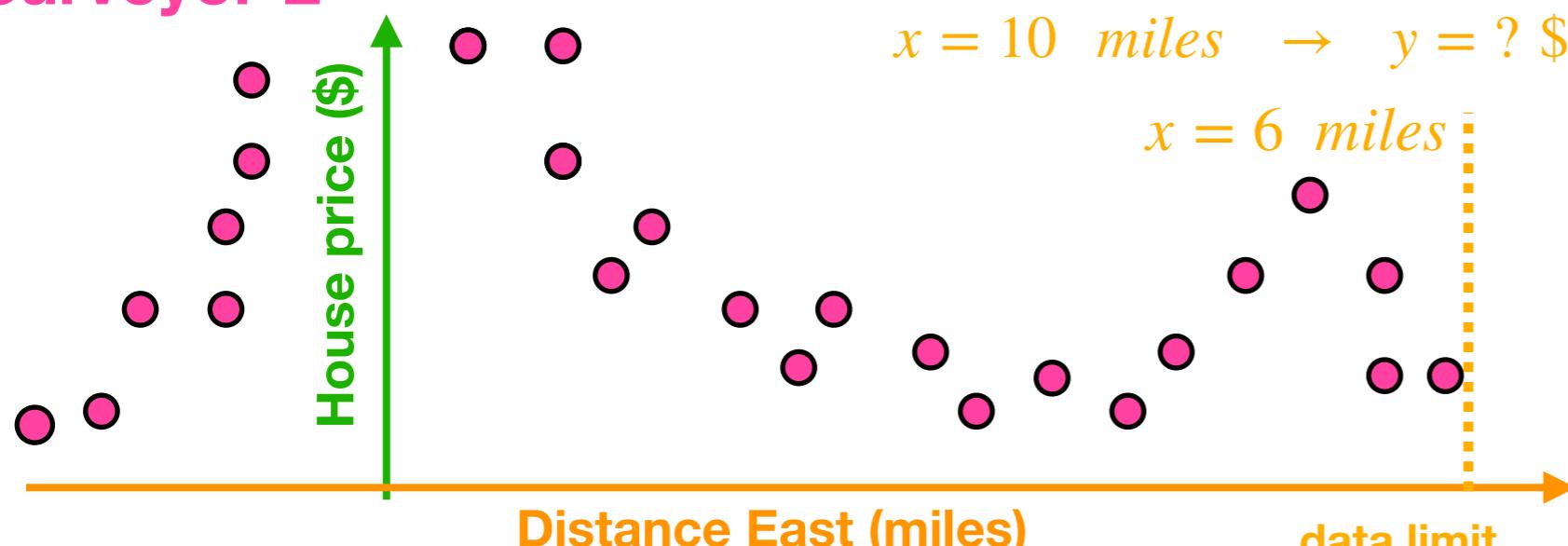
$$\hat{\mathbf{x}} = \begin{pmatrix} 0.25 \\ 0.30 \\ \vdots \\ 3.14 \end{pmatrix} \quad \begin{matrix} \text{Sample-1} \\ \text{Sample-2} \\ \vdots \\ \text{Sample-50} \end{matrix}$$

"Sample" dimension

$$\hat{\mathbf{y}} = \begin{pmatrix} 500,000 \\ 410,000 \\ \vdots \\ 210,000 \end{pmatrix} \quad \begin{matrix} \text{Sample-1} \\ \text{Sample-2} \\ \vdots \\ \text{Sample-50} \end{matrix}$$

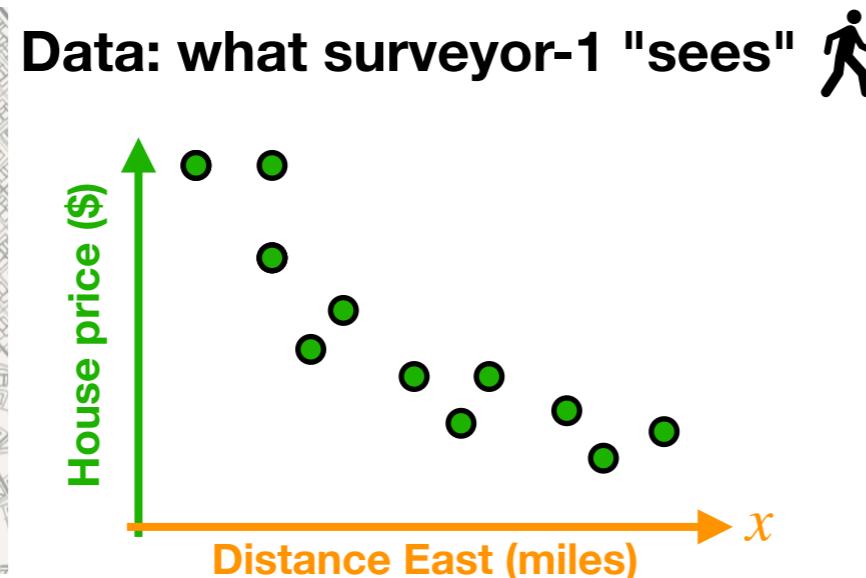
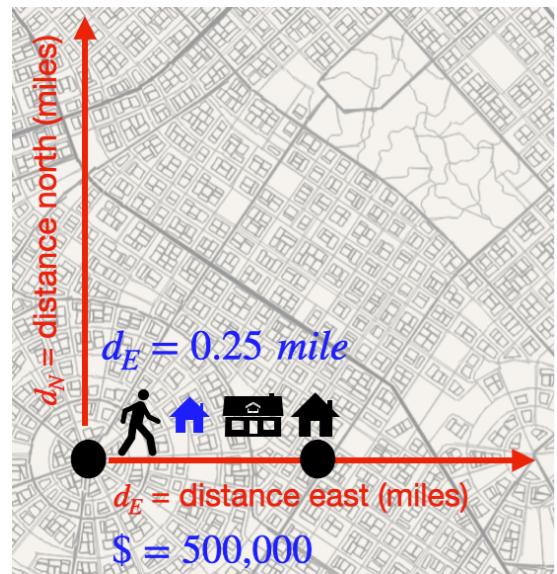


Surveyor-2

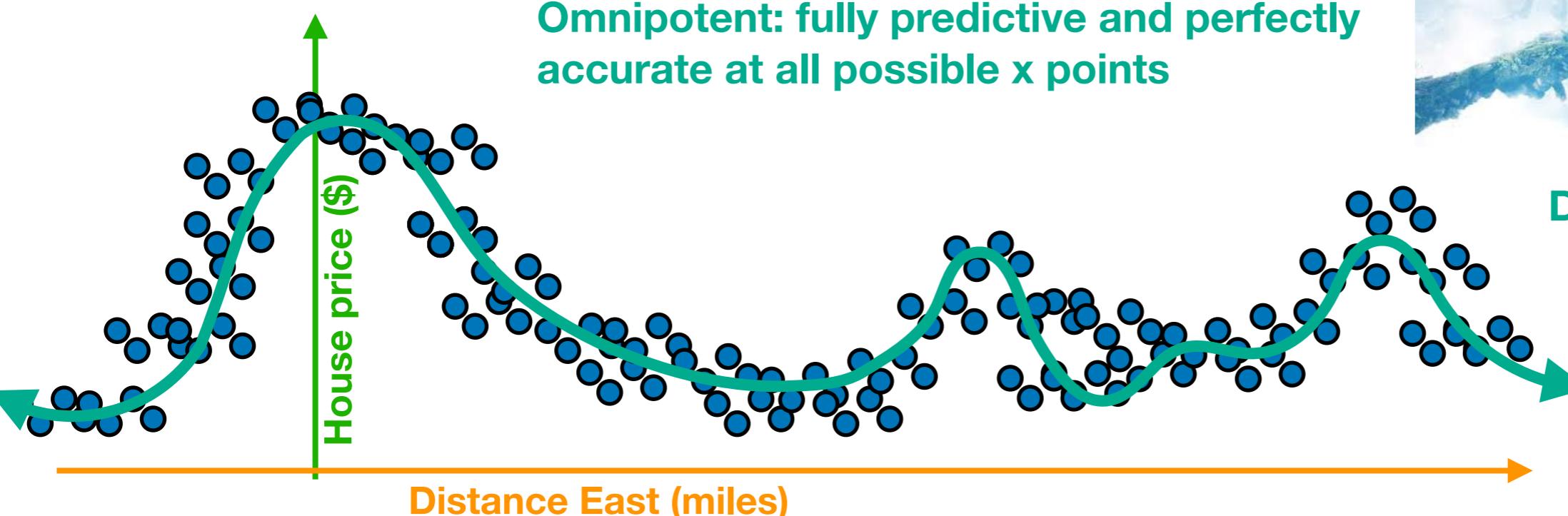


- Data itself generally has little **predictive power**
 - Can only see **local** qualitative trends
 - Are the values big or small
 - Are they going up or down?
 - Are they changing rapidly or slowly?

Data and the “ground truth”



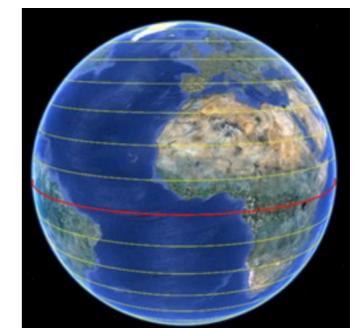
Data generating function: $G(x)$ → UNKNOWN “ground truth”



Omnipotent: fully predictive and perfectly accurate at all possible x points



Defined everywhere

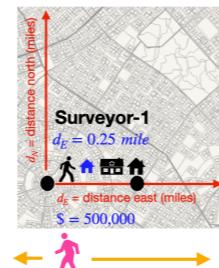


Ideal data → infinite sampling capability $N \rightarrow \infty$

Common ML workflow summary

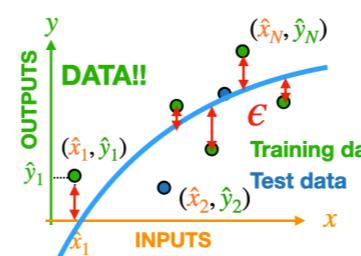
Source and recommended reading: Chollet p111-p116

- Defining the problem and assembling a dataset



- Choosing a measure of success

- Choosing the accuracy/error metrics and a loss function



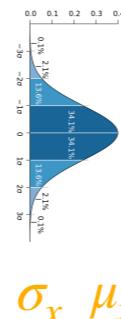
$$\text{RMSE: } \epsilon(\mathbf{p}) = \left(\frac{\sum_i^N (y(\hat{x}_i | \mathbf{p}) - \hat{y}_i)^2}{N} \right)^{\frac{1}{2}}$$

- Deciding on an evaluation protocol



- Preparing your data
- Makes it easier to choose parameter initial conditions (IC)

$$\hat{\mathbf{x}} = \begin{pmatrix} 0.25 \\ 0.30 \\ \vdots \\ 3.14 \end{pmatrix}$$



$$x_i \rightarrow \tilde{x}_i = \frac{x_i - \mu_x}{\sigma_x} \approx [-1.0, 1.0]$$

$$\hat{\mathbf{y}} = \begin{pmatrix} 500,000 \\ 410,000 \\ \vdots \\ 210,000 \end{pmatrix}$$



$$y_i \rightarrow \tilde{y}_i = \frac{y_i - \mu_y}{\sigma_y} \approx [-1.0, 1.0]$$

- Developing a model that does better than some baseline

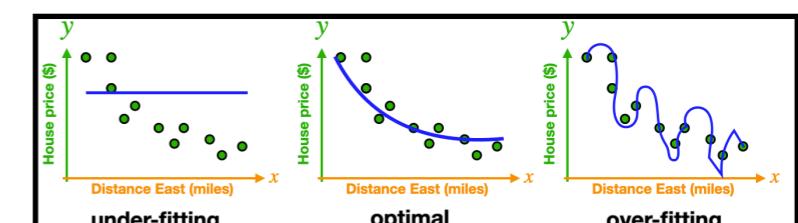
Select a model

$M(x | \mathbf{p})$ = Model = Functional form + Parameterization

$$M(x | \mathbf{p}) = p_1 e^{-\left(\frac{x-p_2}{p_3}\right)^2} + p_4 e^{-\left(\frac{x-p_5}{p_6}\right)^2} + p_7$$

$$\mathbf{p} = (p_0, p_1, p_2, p_3, p_4, p_5, p_6, p_7)$$

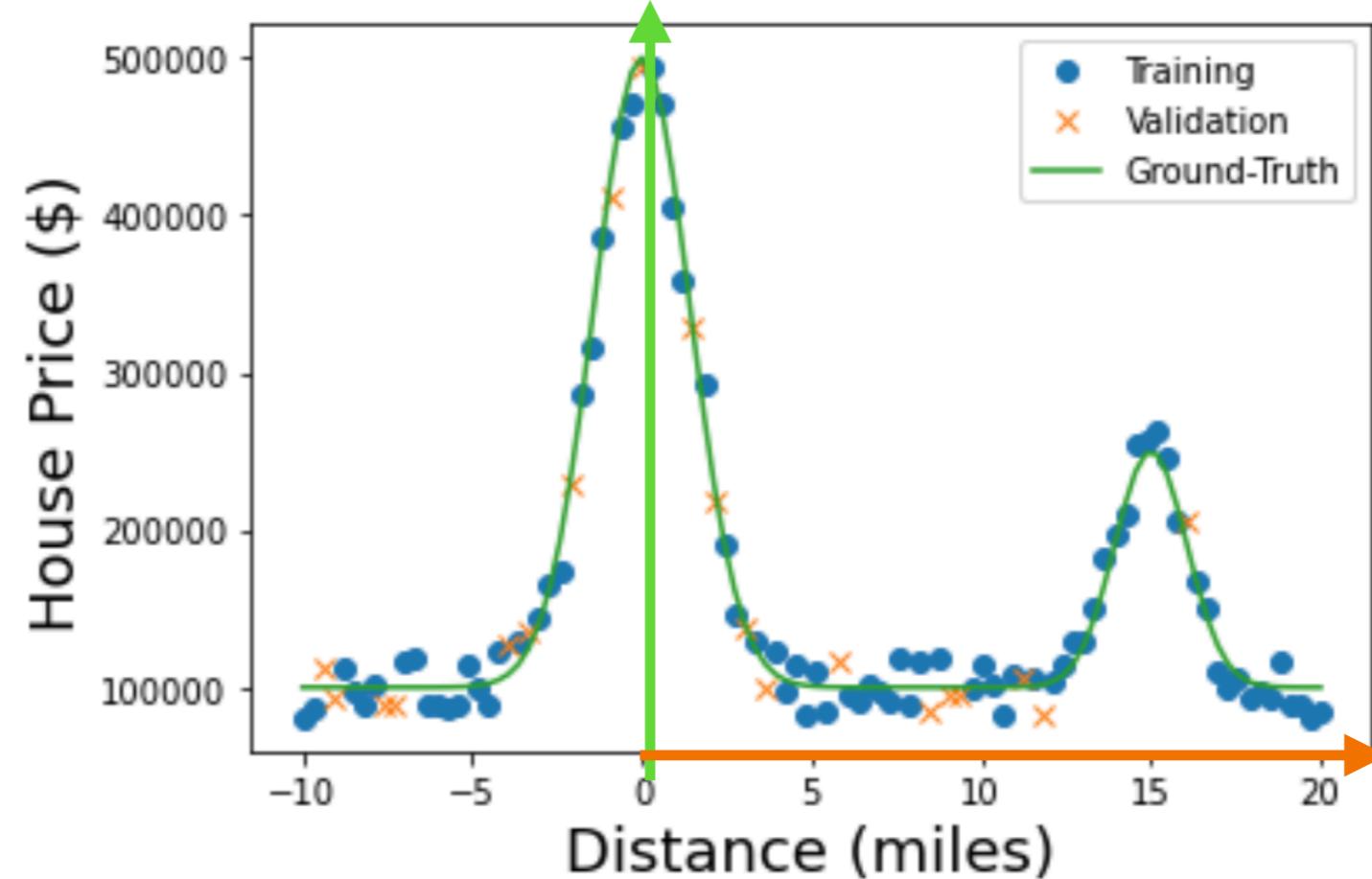
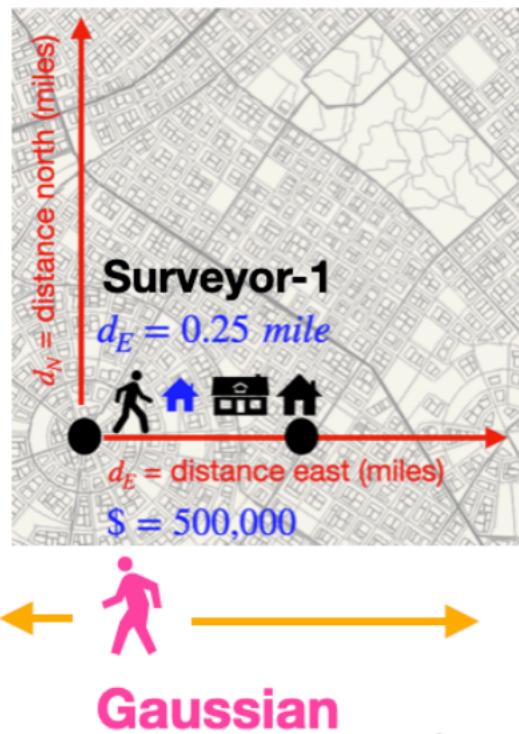
- Choosing an optimizer, learning rate, and other hyper-parameters



★ Scaling up: developing a model that overfits

★ Regularizing your model and tuning your hyper-parameters

Normalize



$$\hat{\mathbf{y}} = \begin{pmatrix} 500,000 \\ 410,000 \\ \vdots \\ 210,000 \end{pmatrix}$$

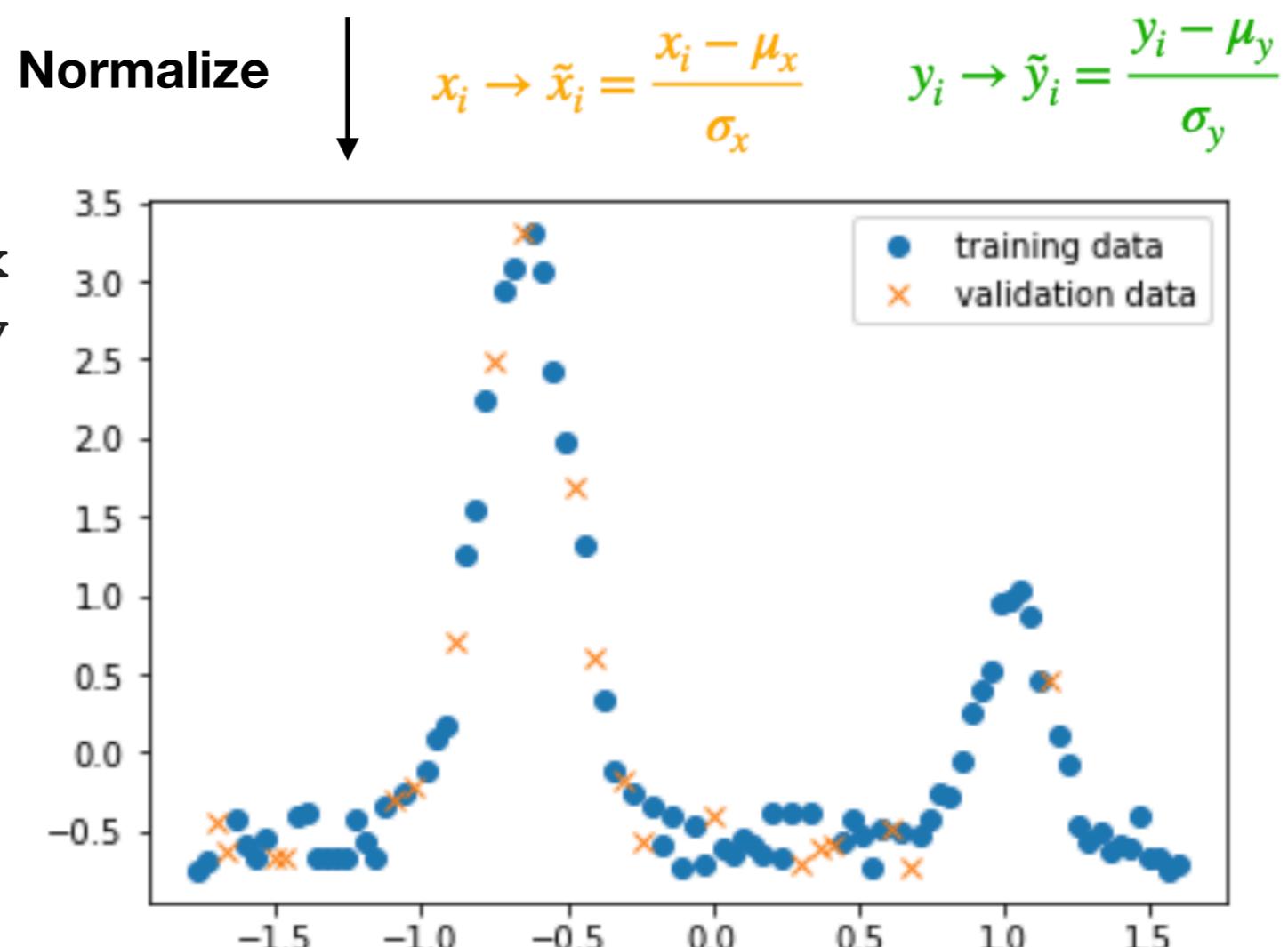
$\sigma_y \quad \mu_y$

$$\hat{\mathbf{x}} = \begin{pmatrix} 0.25 \\ 0.30 \\ \vdots \\ 3.14 \end{pmatrix}$$

$\sigma_x \quad \mu_x$

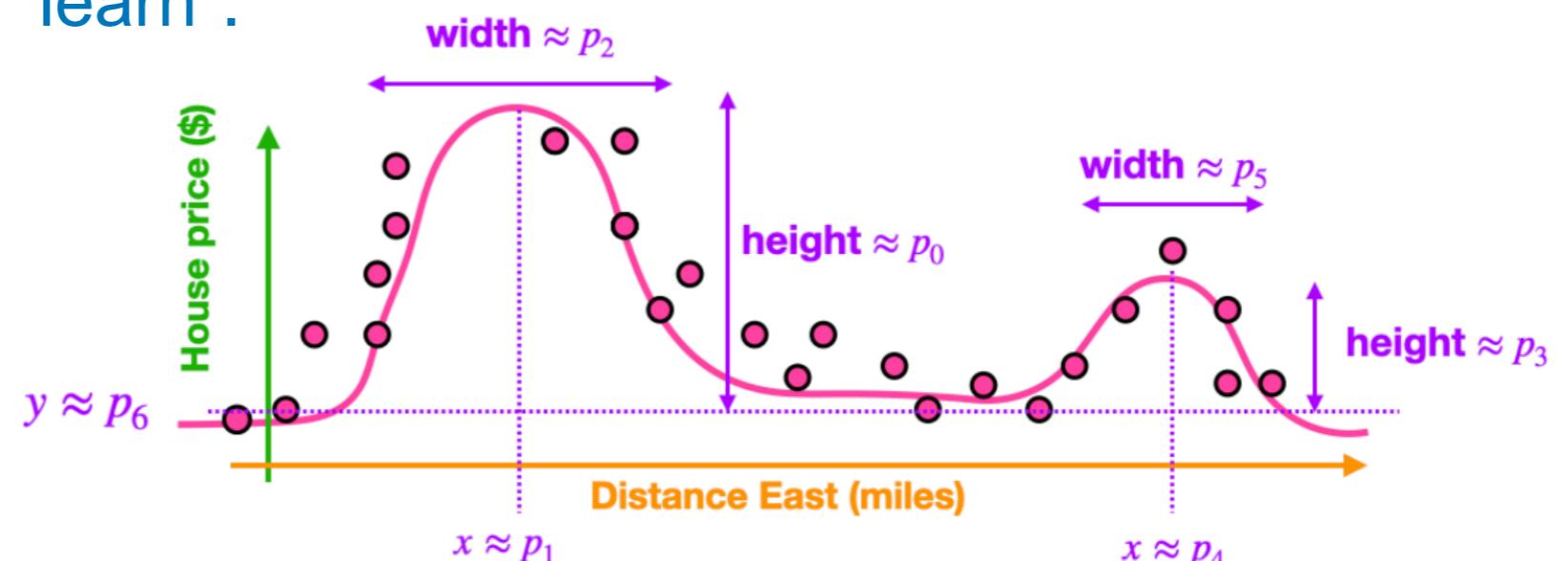
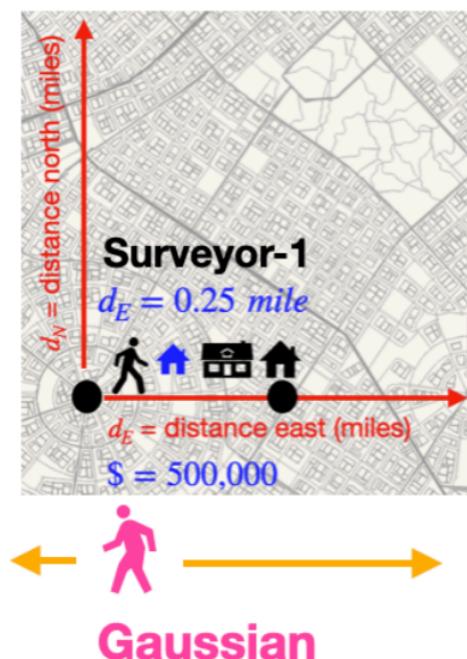
```
ux=np.mean(x_train)
sx=np.std(x_train)
x_train = (x_train-ux)/sx
y_train = (y_train-uy)/sy
```

Transformed that data
Into a dimensionless
space that runs from
roughly -3 to 3 (i.e.
unit-less space)



Why do we normalize?

- 1) Normalization forces the data into a standardized space, centered near the origin (0,0), with a data range \sim (-3 to 3). This effectively “gets rid” of units. However, this range depends on the statistical distribution of the original data
- 2) Normalization makes visualization more effective, since all numerical variables span roughly the same range
- 3) In this normalized space it is MUCH easier for machine learning algorithms to “learn” the data. For example, in parametric modeling we try to “learn” a curve that fits the data by adjusting (i.e. learning) the model parameters. In the normalized space, these parameters will also tend to be in the range (-3 to 3), and therefore easier to “learn”.



$$M(\mathbf{x} | \mathbf{p}) = \text{Model} = \text{Functional form} + \text{Parameterization}$$

$$M(\mathbf{x} | \mathbf{p}) = p_0 e^{-\left(\frac{x-p_1}{p_2}\right)^2} + p_3 e^{-\left(\frac{x-p_4}{p_5}\right)^2} + p_6$$

$$\mathbf{p} = (p_0, p_1, p_2, p_3, p_4, p_5, p_6)$$

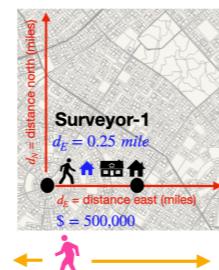
“The goal of normalization is to transform features to be on a similar scale. This improves the performance and training stability of the model”

<https://developers.google.com/machine-learning/data-prep/transform/normalization>

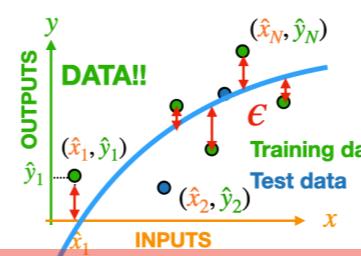
Common ML workflow summary

Source and recommended reading: Chollet p111-p116

- Defining the problem and assembling a dataset

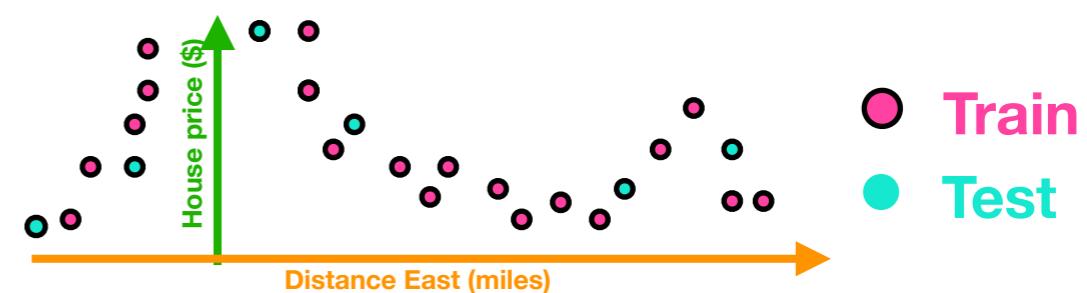


- Choosing a measure of success
 - Choosing the accuracy/error metrics and a loss function



$$\text{RMSE: } \epsilon(\mathbf{p}) = \left(\frac{\sum_i^N (y(\hat{x}_i | \mathbf{p}) - \hat{y}_i)^2}{N} \right)^{\frac{1}{2}}$$

- Deciding on an evaluation protocol



- Preparing your data

Makes it easier to choose parameter initial conditions (IC)

$$\hat{\mathbf{x}} = \begin{pmatrix} 0.25 \\ 0.30 \\ \vdots \\ 3.14 \end{pmatrix} \quad \begin{matrix} \sigma_x & \mu_x \end{matrix} \quad x_i \rightarrow \tilde{x}_i = \frac{x_i - \mu_x}{\sigma_x} \quad \hat{\mathbf{y}} = \begin{pmatrix} 500,000 \\ 410,000 \\ \vdots \\ 210,000 \end{pmatrix} \quad \begin{matrix} \sigma_y & \mu_y \end{matrix}$$

$$y_i \rightarrow \tilde{y}_i = \frac{y_i - \mu_y}{\sigma_y} \approx [-1.0, 1.0]$$

- Developing a model that does better than some baseline

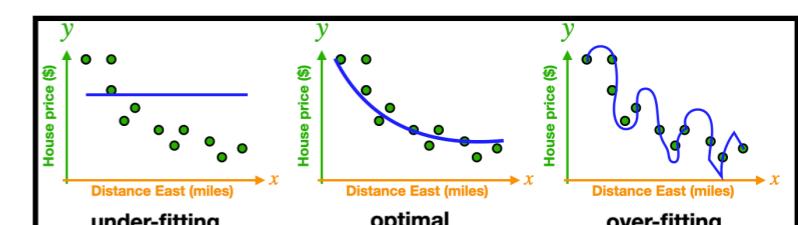
Select a model

$M(x | \mathbf{p})$ = Model = Functional form + Parameterization

$$M(x | \mathbf{p}) = p_1 e^{-\left(\frac{x-p_2}{p_3}\right)^2} + p_4 e^{-\left(\frac{x-p_5}{p_6}\right)^2} + p_7$$

$$\mathbf{p} = (p_0, p_1, p_2, p_3, p_4, p_5, p_6, p_7)$$

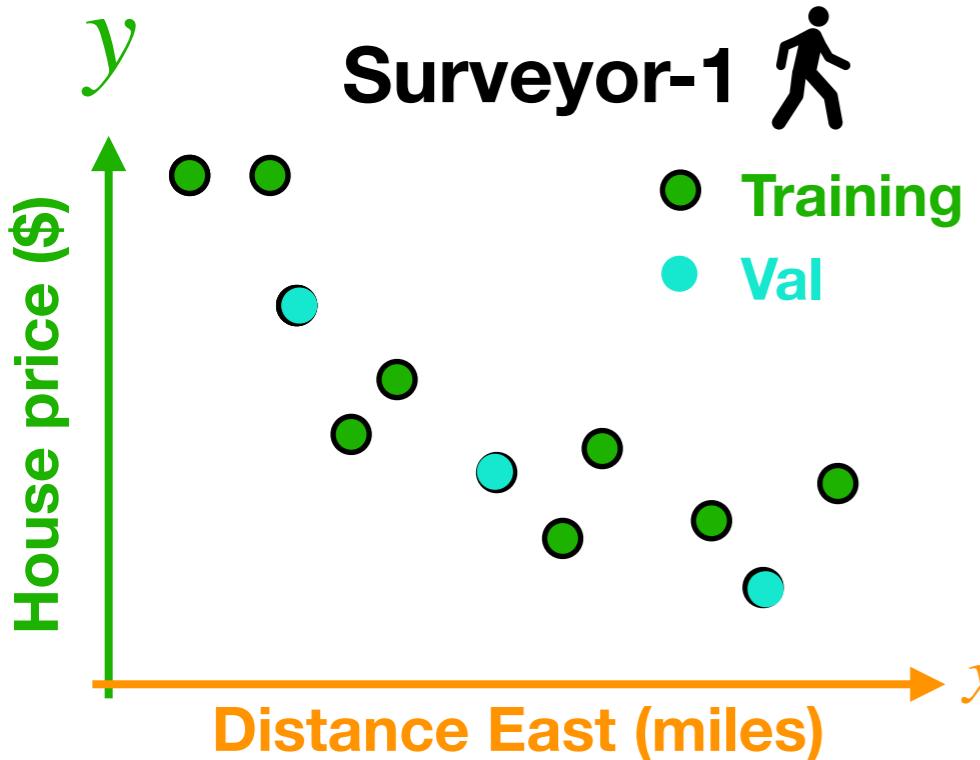
- Choosing an optimizer, learning rate, and other hyper-parameters



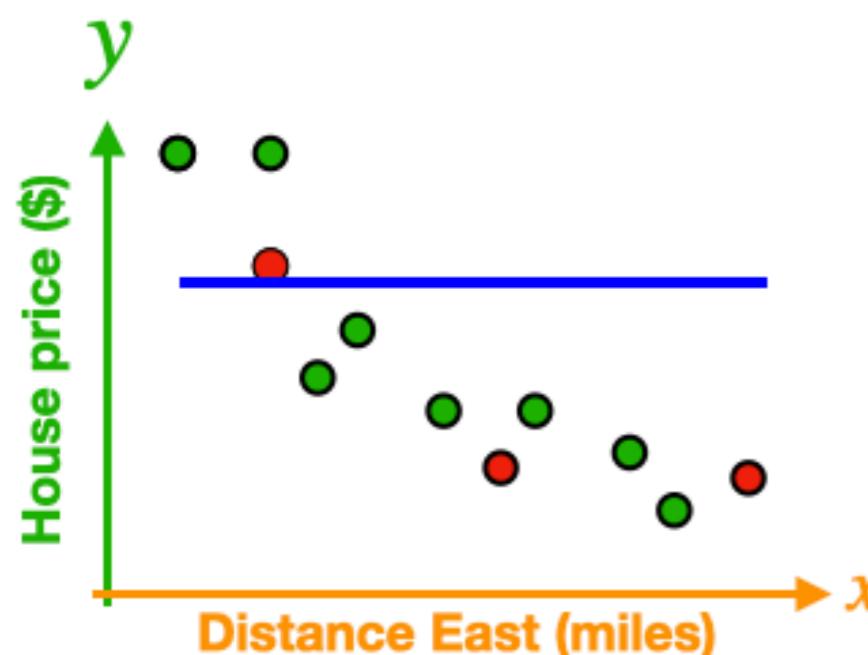
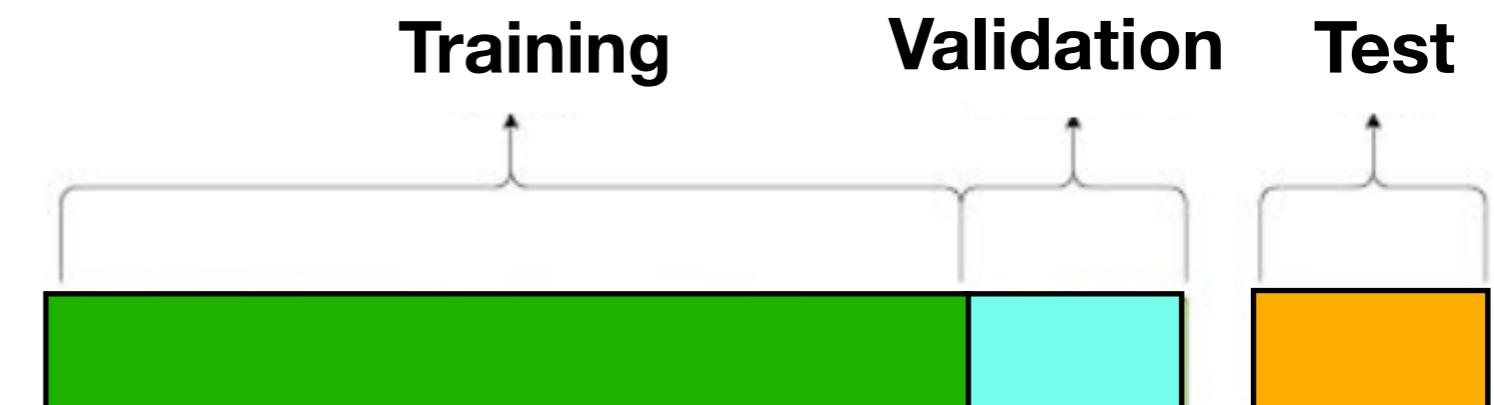
★ Scaling up: developing a model that overfits

★ Regularizing your model and tuning your hyper-parameters

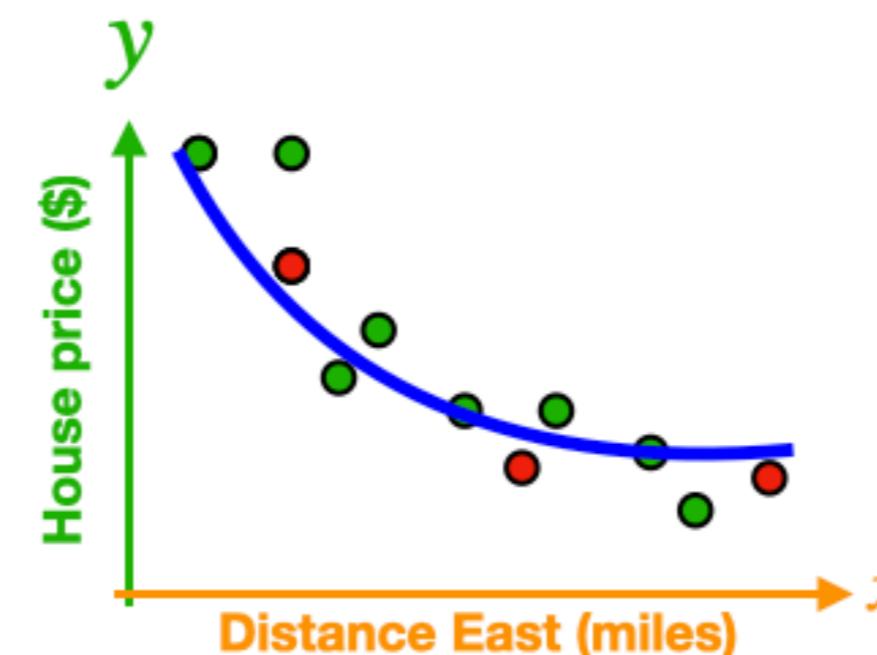
Partition your data



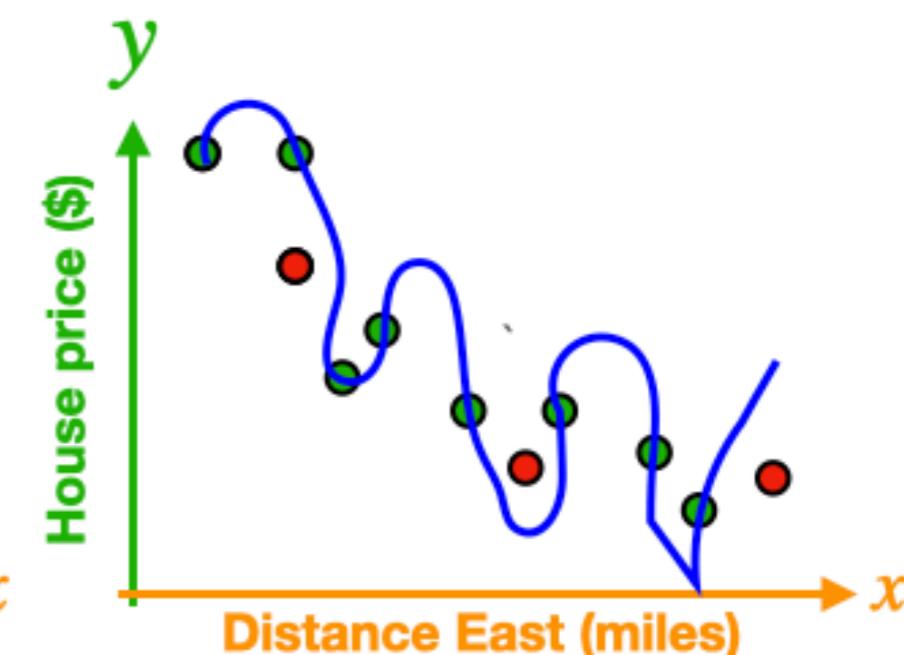
1) Partition data



under-fitting
model complexity is
too low



optimal
correct level of
model complexity

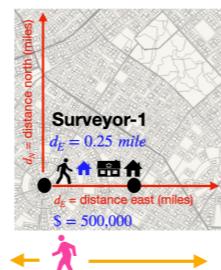


over-fitting
model complexity
is too high

Common ML workflow summary

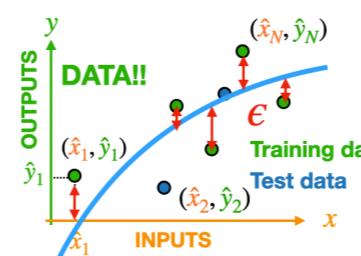
Source and recommended reading: Chollet p111-p116

- Defining the problem and assembling a dataset



- Choosing a measure of success

- Choosing the accuracy/error metrics and a loss function



$$\text{RMSE: } \epsilon(\mathbf{p}) = \left(\frac{\sum_i^N (y(\hat{x}_i | \mathbf{p}) - \hat{y}_i)^2}{N} \right)^{\frac{1}{2}}$$

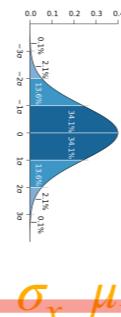
- Deciding on an evaluation protocol



- Preparing your data

Makes it easier to choose parameter initial conditions (IC)

$$\hat{\mathbf{x}} = \begin{pmatrix} 0.25 \\ 0.30 \\ \vdots \\ 3.14 \end{pmatrix}$$



$$x_i \rightarrow \tilde{x}_i = \frac{x_i - \mu_x}{\sigma_x}$$

$$\approx [-1.0, 1.0]$$

$$\hat{\mathbf{y}} = \begin{pmatrix} 500,000 \\ 410,000 \\ \vdots \\ 210,000 \end{pmatrix}$$

$$y_i \rightarrow \tilde{y}_i = \frac{y_i - \mu_y}{\sigma_y}$$

$$\approx [-1.0, 1.0]$$

- Developing a model that does better than some baseline

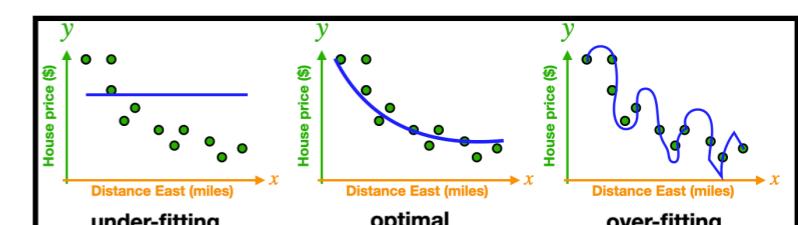
Select a model

$M(x | \mathbf{p}) = \text{Model} = \text{Functional form} + \text{Parameterization}$

$$M(x | \mathbf{p}) = p_1 e^{-\left(\frac{x-p_2}{p_3}\right)^2} + p_4 e^{-\left(\frac{x-p_5}{p_6}\right)^2} + p_7$$

$$\mathbf{p} = (p_0, p_1, p_2, p_3, p_4, p_5, p_6, p_7)$$

- Choosing an optimizer, learning rate, and other hyper-parameters



★ Scaling up: developing a model that overfits

★ Regularizing your model and tuning your hyper-parameters

Choosing a model

$$\text{MODEL} = \text{Functional form} + \text{Parameterization} \quad y = M(x | p)$$

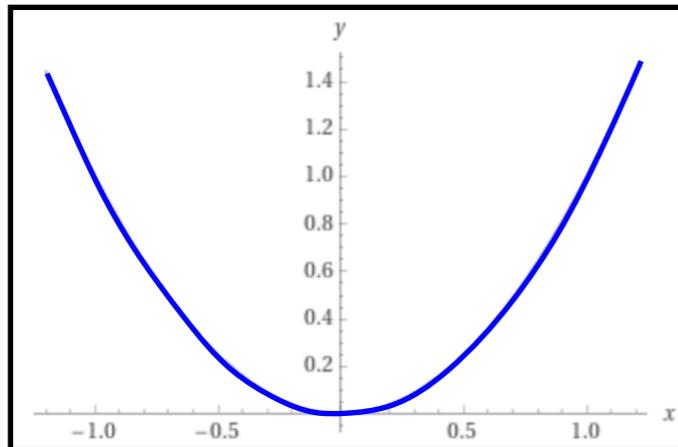
(Parametric)

(underlying shape)

(stretching and shifting)

Surveyor-1

(quadratic)

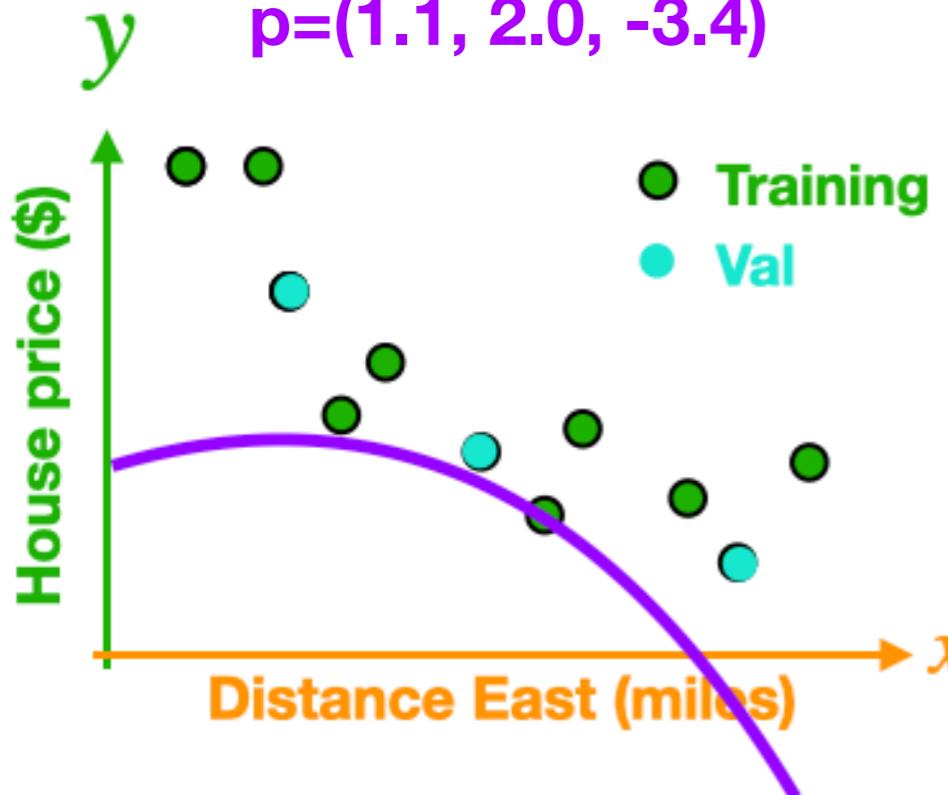


$$M(x | p) = p_0 + p_1 x + p_2 x^2$$
$$p = (p_0, p_1, p_2)$$

We have three parameters, i.e.
“knobs”, to “turn” to fit the data
(i.e degrees of freedom)

Bad parameterization

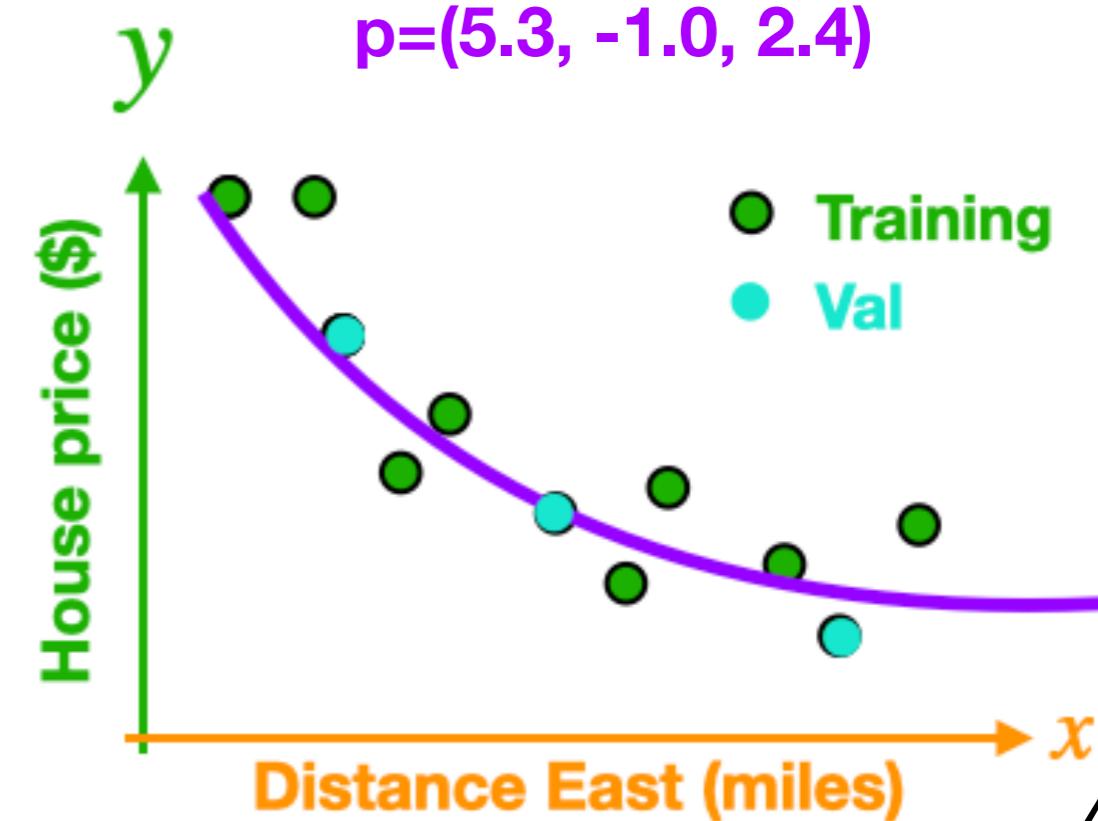
$$p = (1.1, 2.0, -3.4)$$



Goal: Find the three PARAMETER numbers $p = (p_0, p_1, p_2)$ which provide the “best” fit of the data

Good parameterization

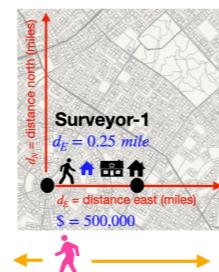
$$p = (5.3, -1.0, 2.4)$$



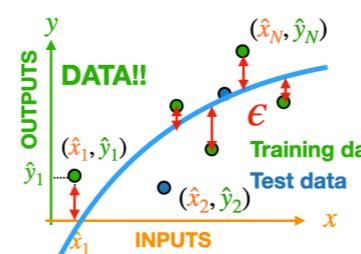
Common ML workflow summary

Source and recommended reading: Chollet p111-p116

- Defining the problem and assembling a dataset



- Choosing a measure of success
 - Choosing the accuracy/error metrics and a loss function



$$\text{RMSE: } \epsilon(\mathbf{p}) = \left(\frac{\sum_i^N (y(\hat{x}_i | \mathbf{p}) - \hat{y}_i)^2}{N} \right)^{\frac{1}{2}}$$

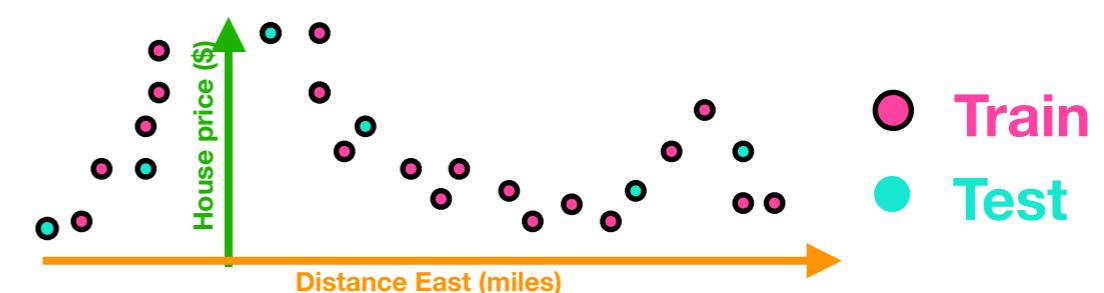
- Deciding on an evaluation protocol



- Preparing your data

Makes it easier to choose parameter initial conditions (IC)

$$\hat{x} = \begin{pmatrix} 0.25 \\ 0.30 \\ \vdots \\ 3.14 \end{pmatrix} \quad \sigma_x \quad \mu_x$$



Train
Test

$$x_i \rightarrow \tilde{x}_i = \frac{x_i - \mu_x}{\sigma_x} \approx [-1.0, 1.0] \quad \hat{y} = \begin{pmatrix} 500,000 \\ 410,000 \\ \vdots \\ 210,000 \end{pmatrix} \quad \sigma_y \quad \mu_y$$

$$y_i \rightarrow \tilde{y}_i = \frac{y_i - \mu_y}{\sigma_y} \approx [-1.0, 1.0]$$

- Developing a model that does better than some baseline

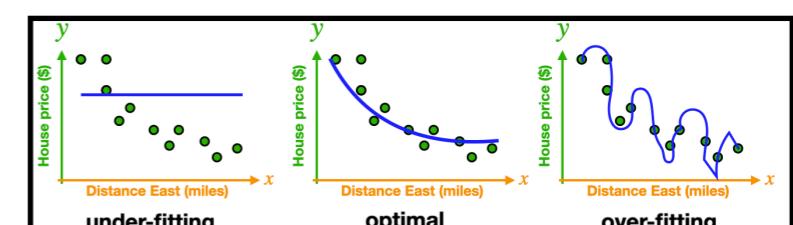
Select a model

$M(x | \mathbf{p})$ = Model = Functional form + Parameterization

$$M(x | \mathbf{p}) = p_1 e^{-\left(\frac{x-p_2}{p_3}\right)^2} + p_4 e^{-\left(\frac{x-p_5}{p_6}\right)^2} + p_7$$

$$\mathbf{p} = (p_0, p_1, p_2, p_3, p_4, p_5, p_6, p_7)$$

- Choosing an optimizer, learning rate, and other hyper-parameters

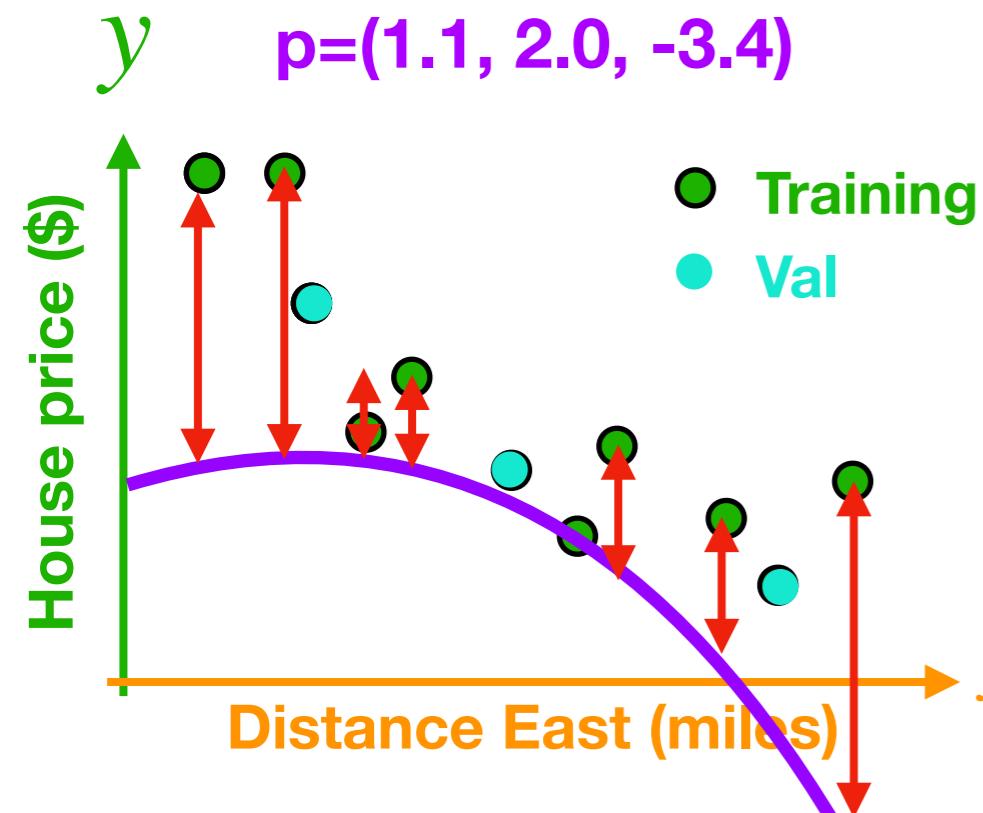


- Scaling up: developing a model that overfits
- Regularizing your model and tuning your hyper-parameters

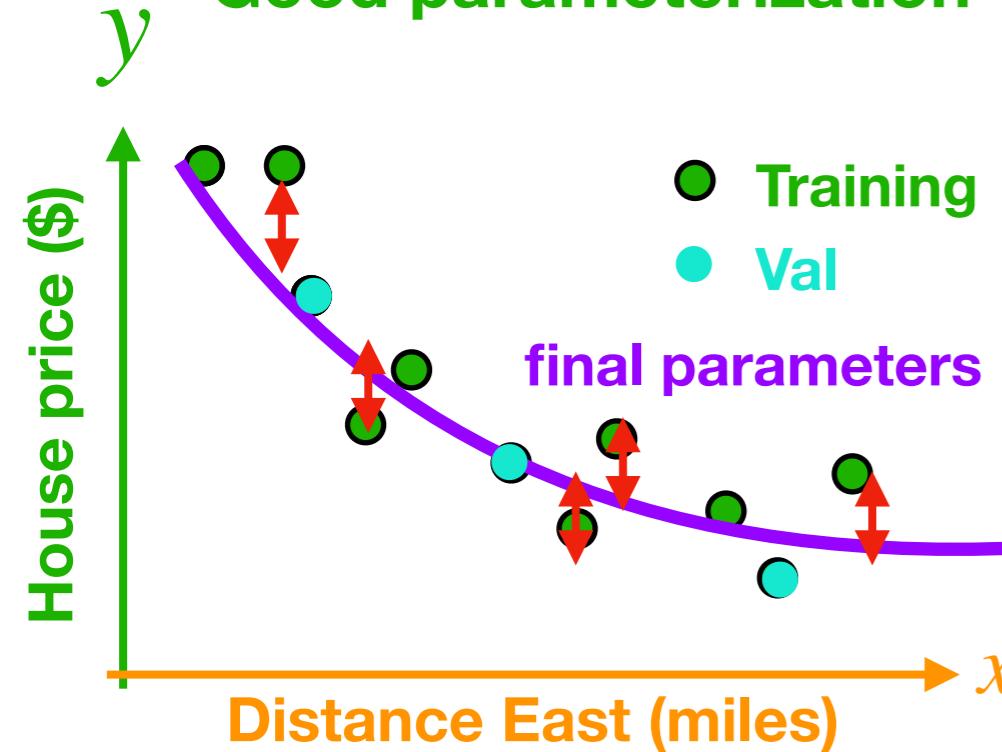
Loss function

Model: $M(x | p) = p_0 + p_1x + p_2x^2$

Bad parameterization

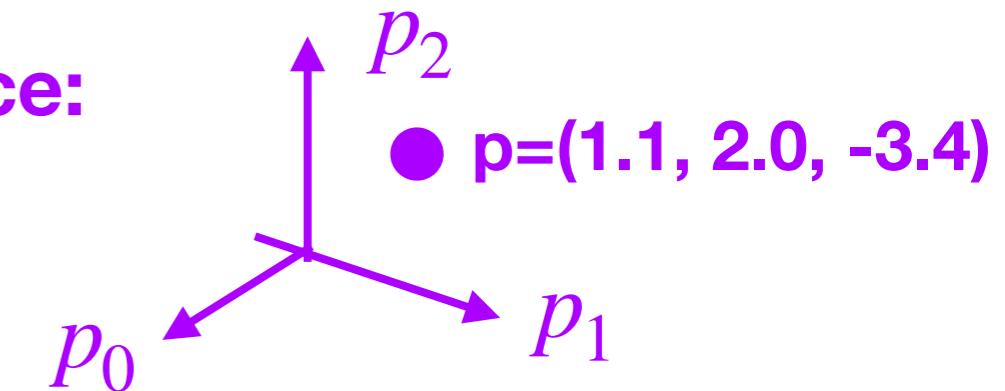


Good parameterization



Parameter space:

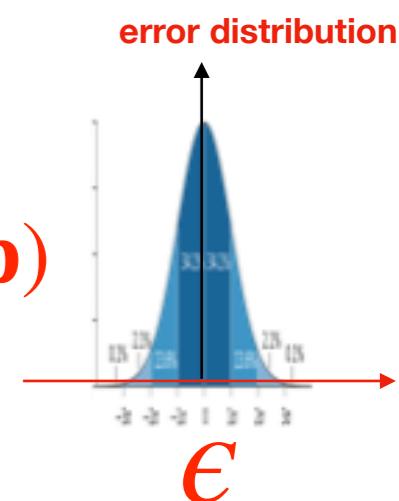
$$p = (p_0, p_1, p_2)$$



Space of all possible quadratic models

Vector of errors (residuals)

$$\epsilon_i(p) = y_{data} - y_{pred} = y_i - M(x_i | p)$$



Define an SCALAR error metric (loss function)

Loss function:

$$L_{train}(p)$$

RMSE:

$$L(p) = \sqrt{\frac{\sum_i^N M(x | p) - \hat{y}_i)^2}{N}}$$

AKA: Objective function

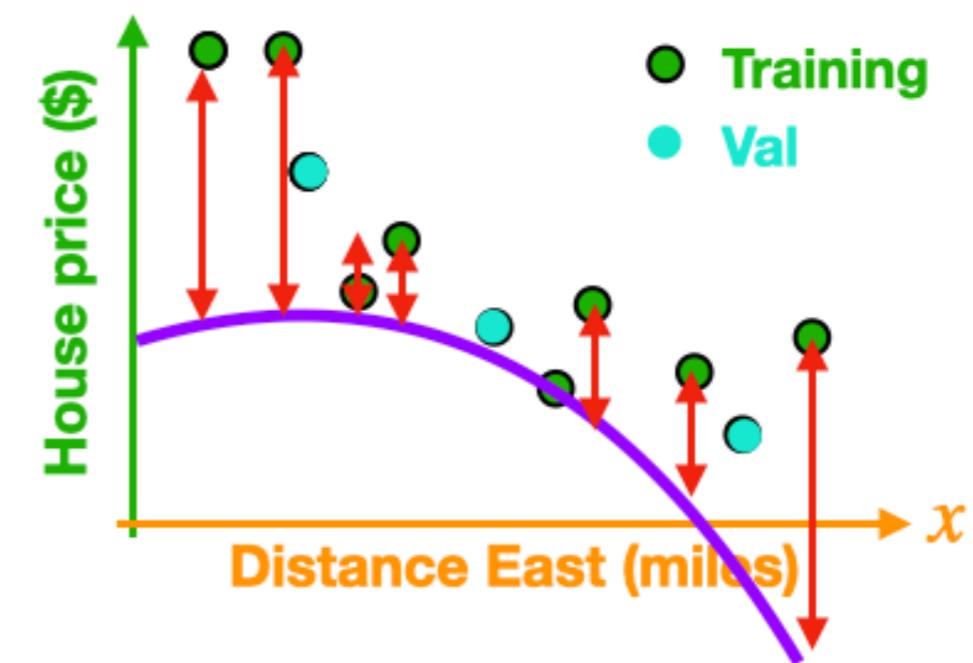
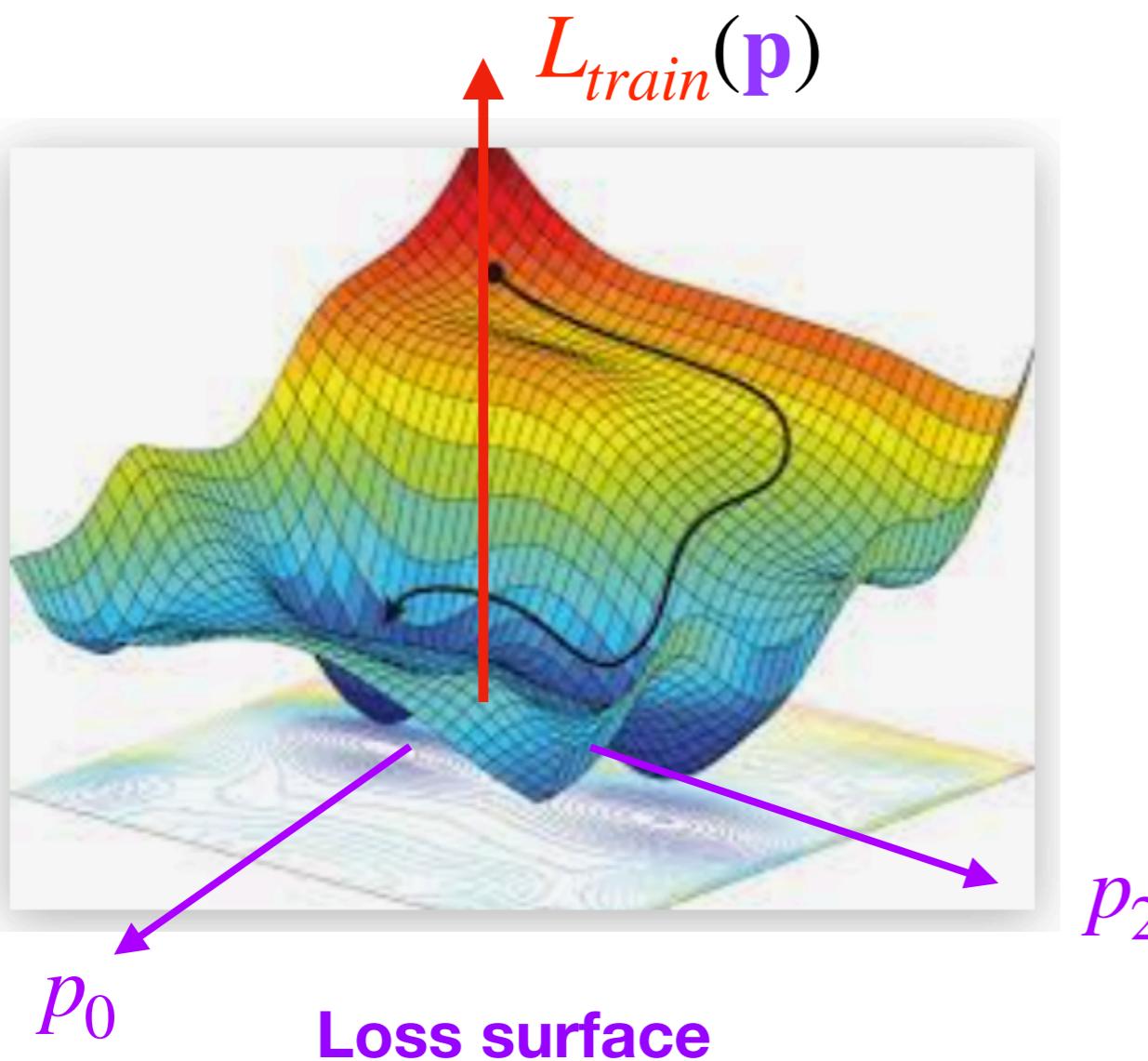
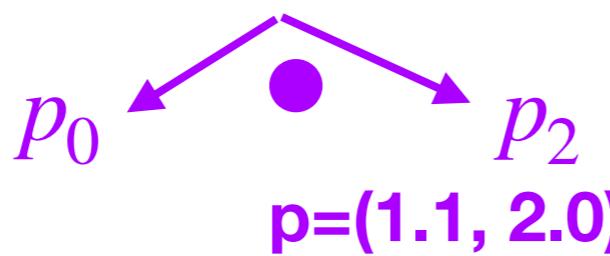
Loss surface-1

simplified model model:

$$M(x | p) = p_0 + p_2 x^2$$

Parameter space:

$$p = (p_0, p_1)$$



$$\epsilon_i(p) = y_{data} - y_{pred} = y_i - M(x_i | p)$$

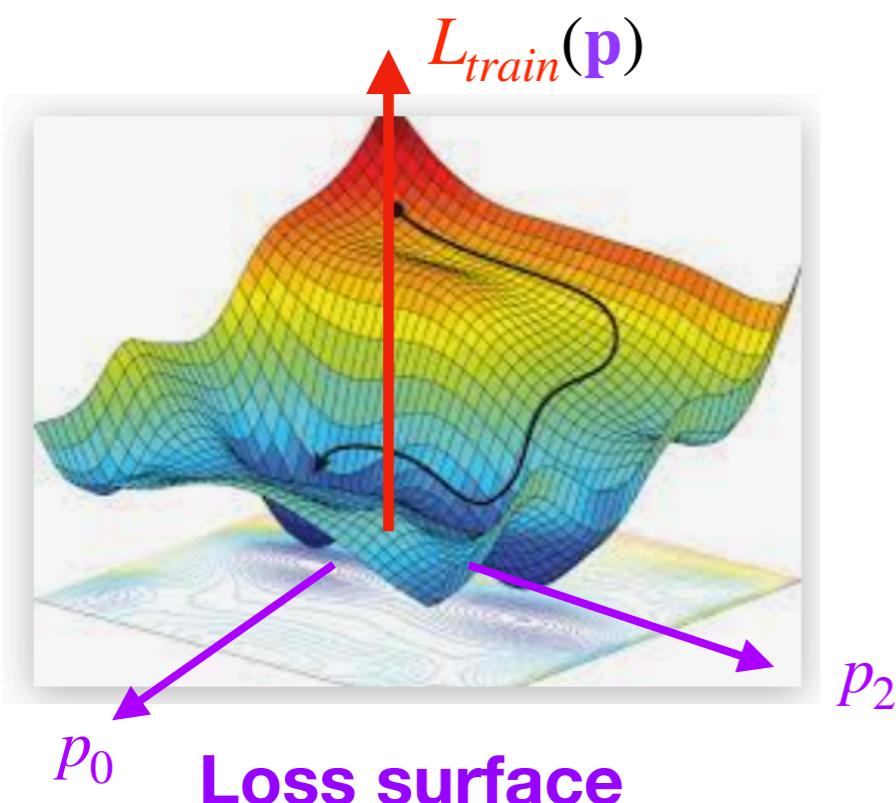
Loss function:

RMSE:

$$L(p) = \left(\frac{\sum_i^N M(x | p) - \hat{y}_i)^2}{N} \right)^{\frac{1}{2}}$$

AKA: Objective function

Training Numerical optimization problem to find best parameters (best fit)



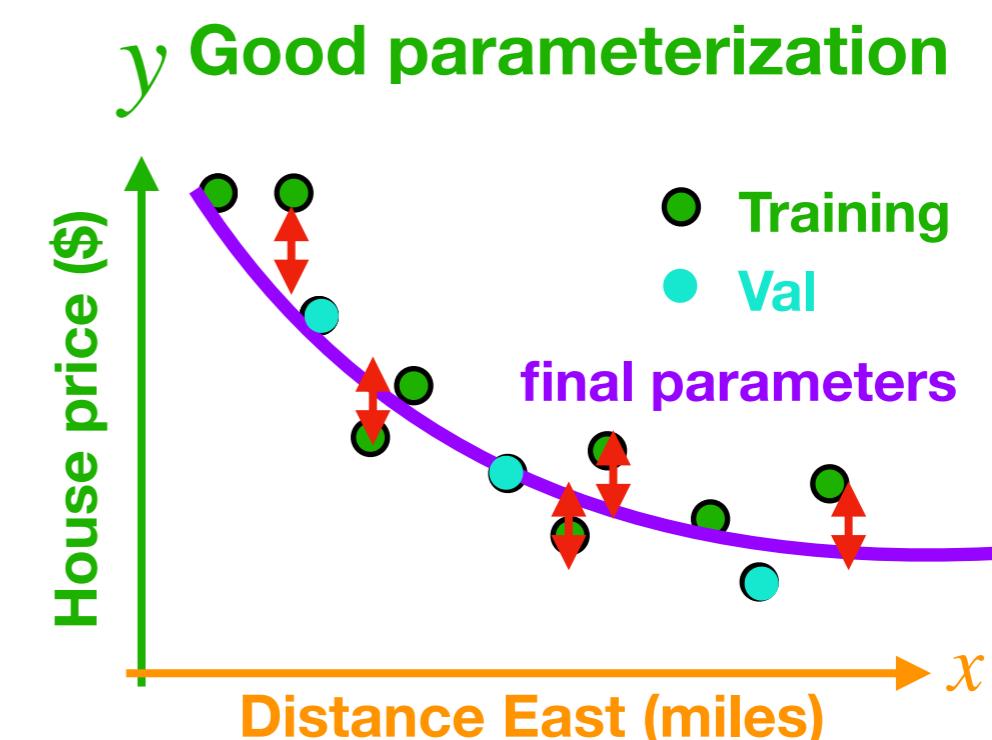
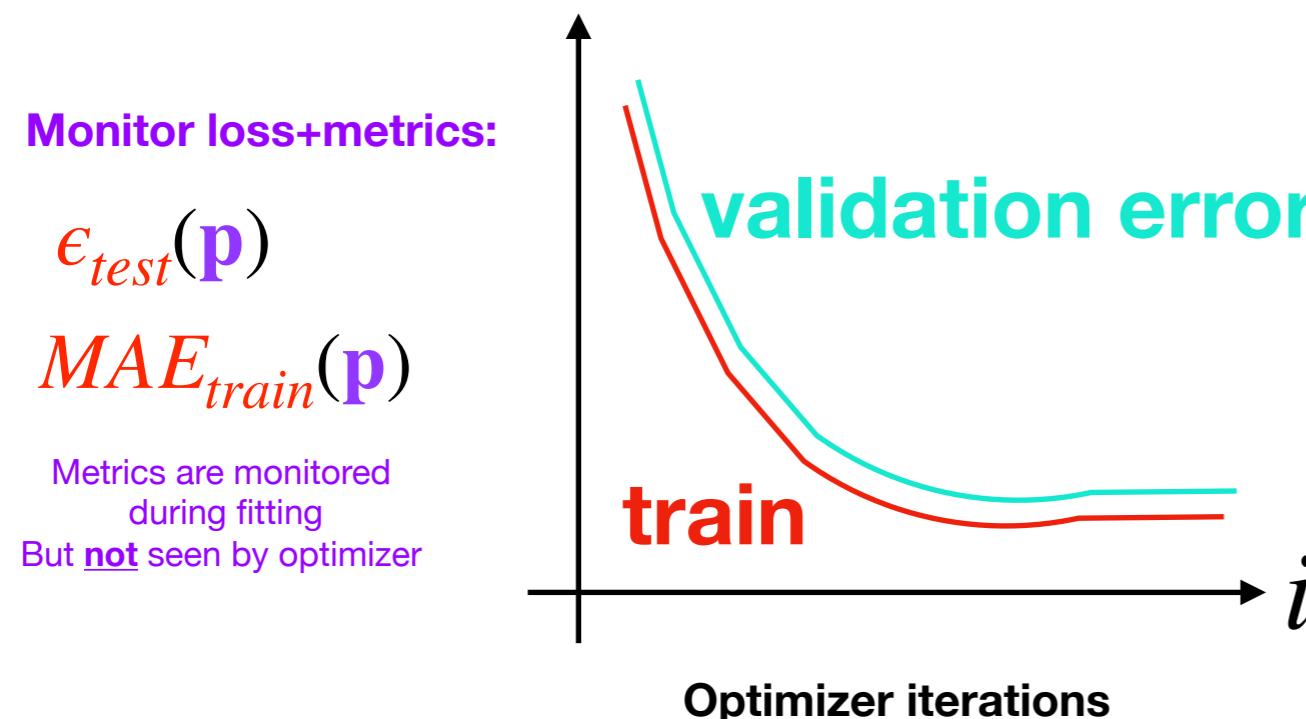
simplified model model:

$$M(x | p) = p_0 + p_2 x^2$$

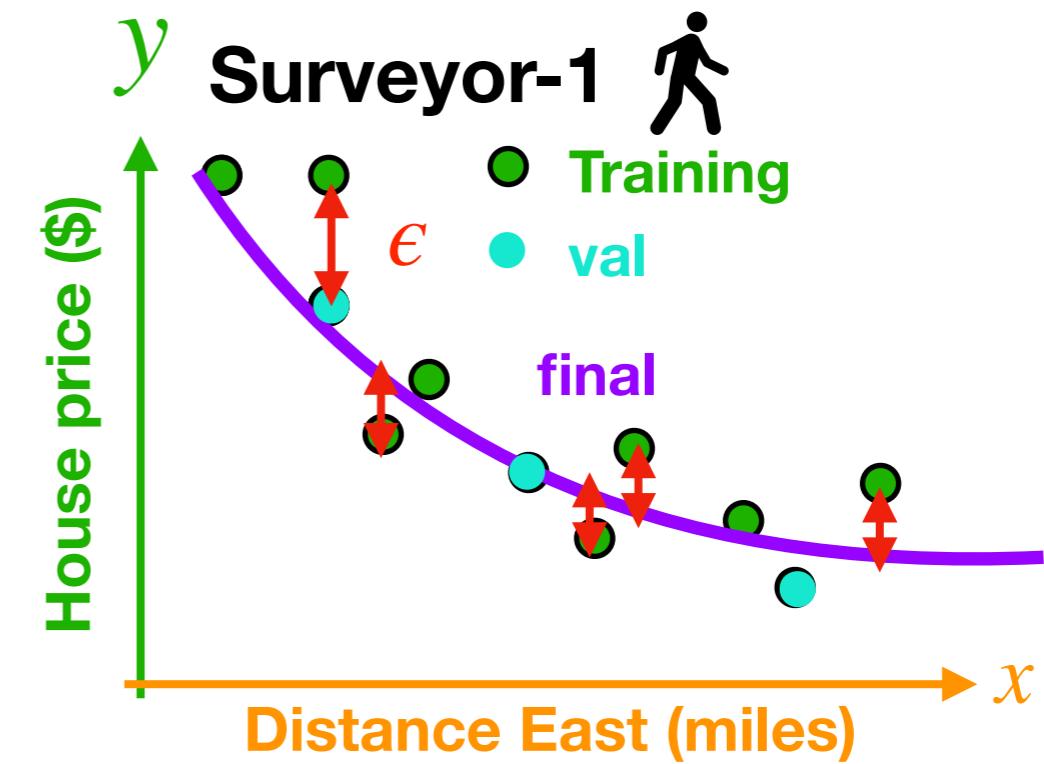
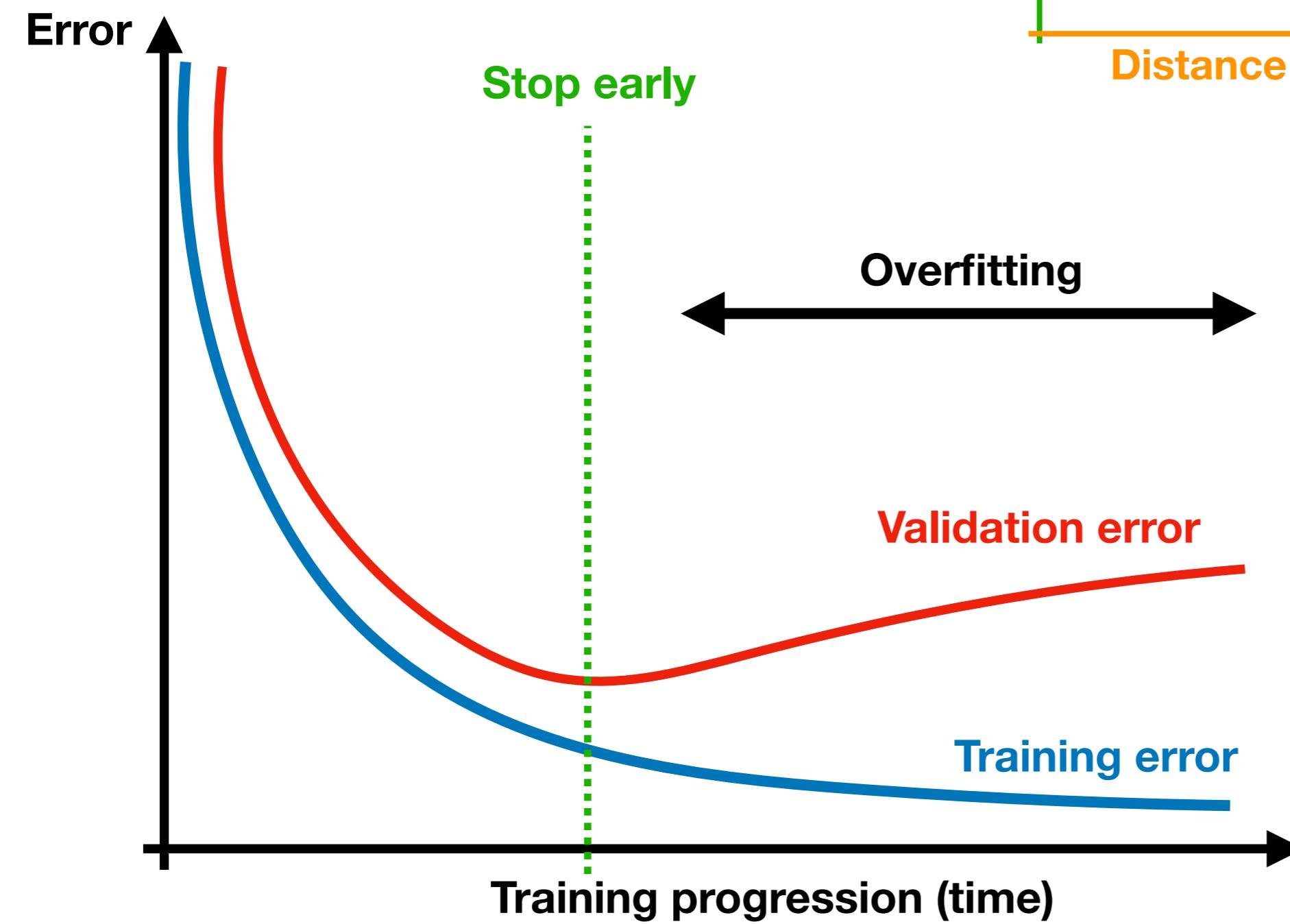
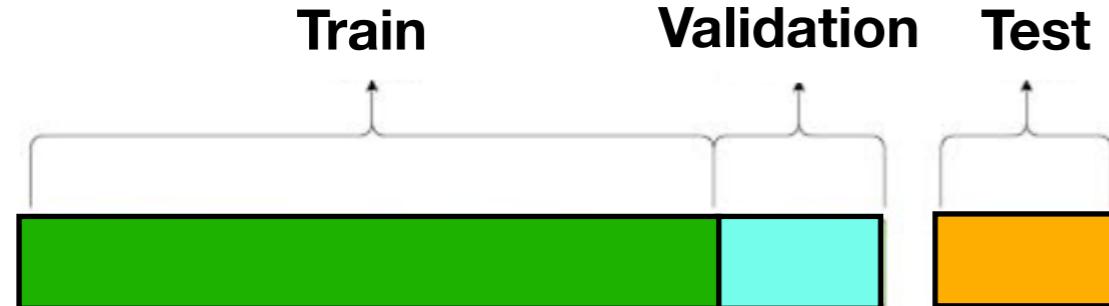
Parameter space: $p = (p_0, p_1)$

Goal find parameters that are local or global minima of the loss surface. The lower the minima the better the model fit

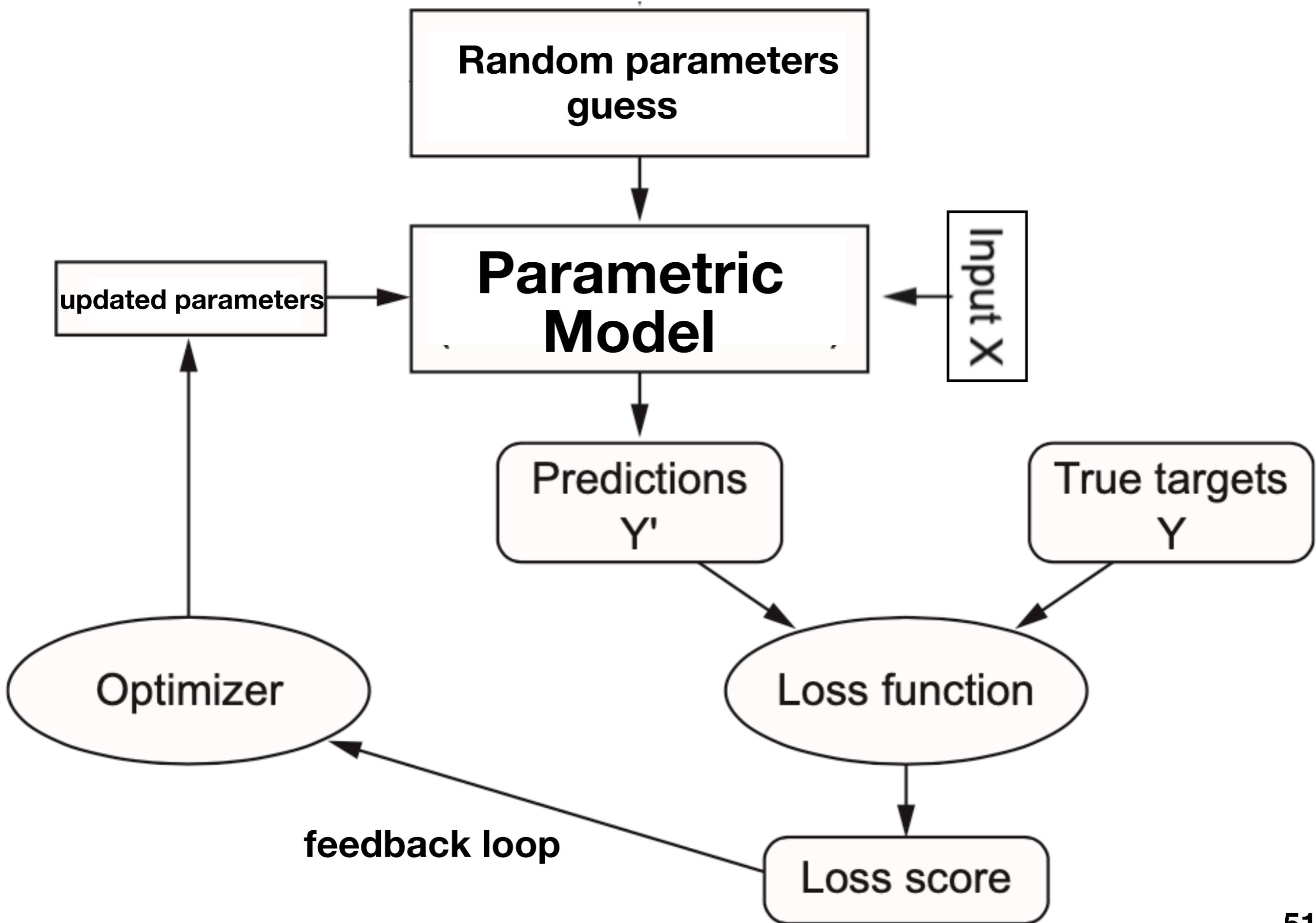
Minimize loss surface with gradient decent $\mathbf{p}_{n+1} = \mathbf{p}_n - \lambda \nabla_w \epsilon(\mathbf{p}_n)$ $\nabla_p \epsilon = \left(\frac{\partial \epsilon}{\partial p_0}, \frac{\partial \epsilon}{\partial p_2} \right)$



Monitoring overfitting

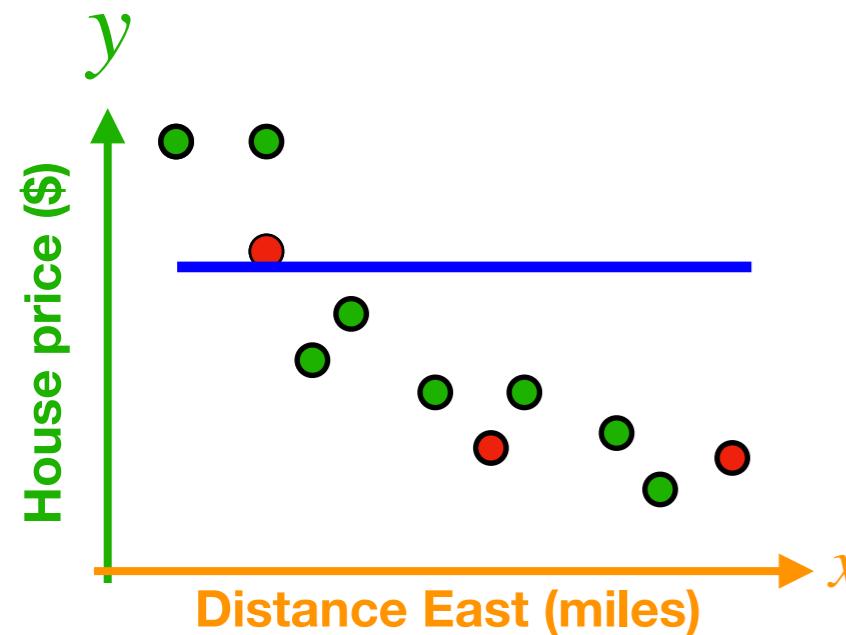


Model training diagram



OVERFITTING

Interpolation concerns



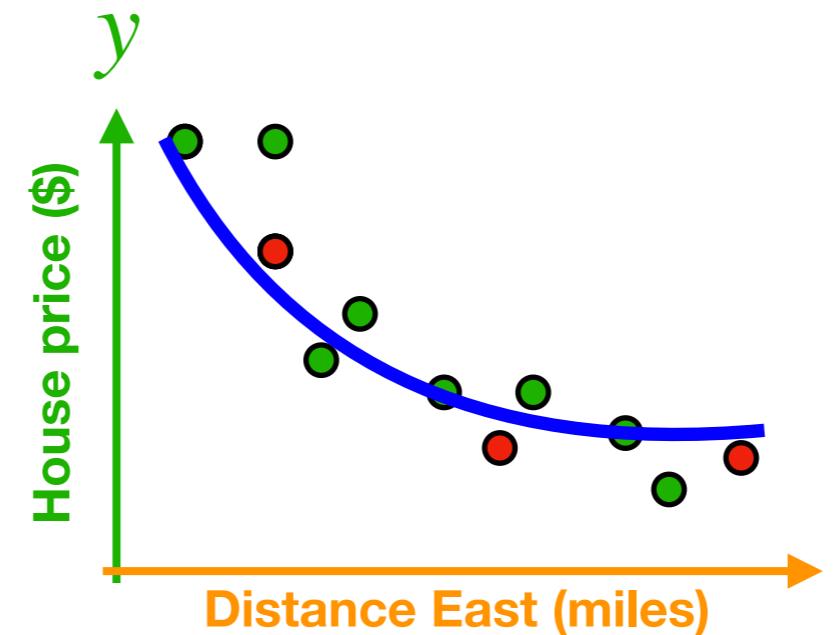
under-fitting
model complexity is
too low

High training error
High test error

high bias
low variance

The inability for a machine learning method (like linear regression) to capture the true relationship is called **bias**.

(Under parameterized)



optimal
correct level of
model complexity

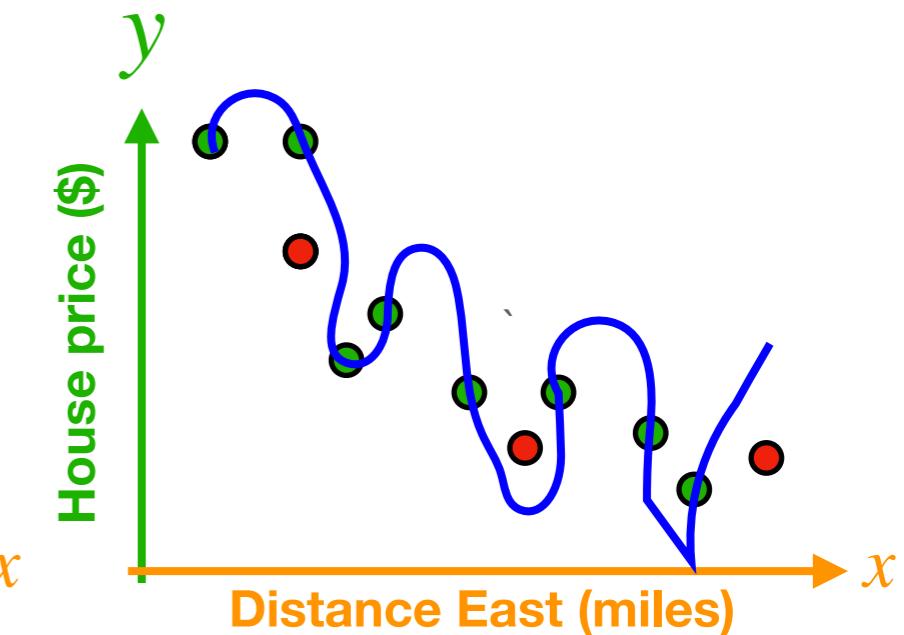
Training and Test error
are equally low

low bias
low variance

In machine learning, the ideal algorithm has **low bias** and can accurately model the true relationship...

...and it has **low variability**, by producing consistent predictions across different datasets.

(Properly parameterized)



over-fitting
model complexity
is too high

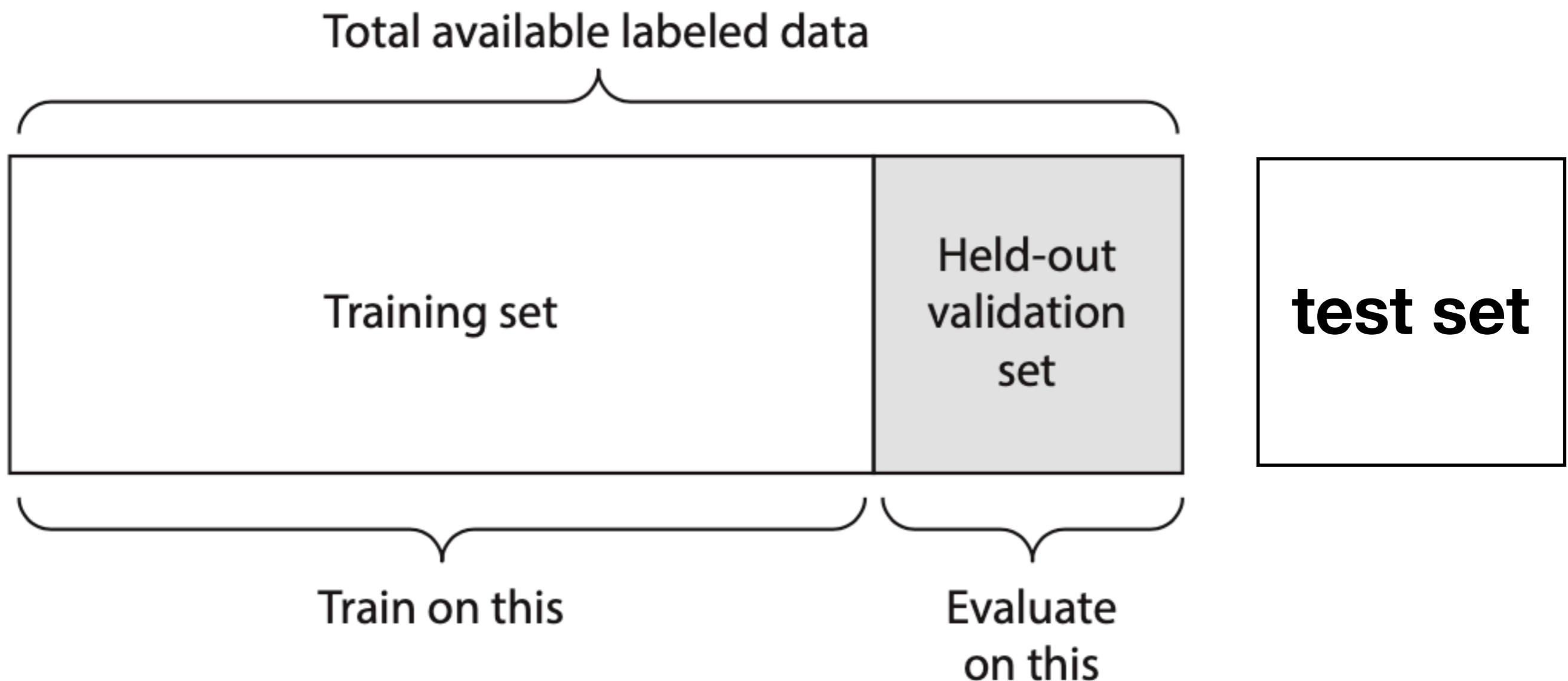
Very Low training error
Very high test error

low bias
high variance

In Machine Learning lingo, the difference in fits between data sets is called **Variance**.

(Over parameterized)

Hold out validation



Each is a random subset of total dataset

General regression concerns

Extrapolation: Predictions "outside"

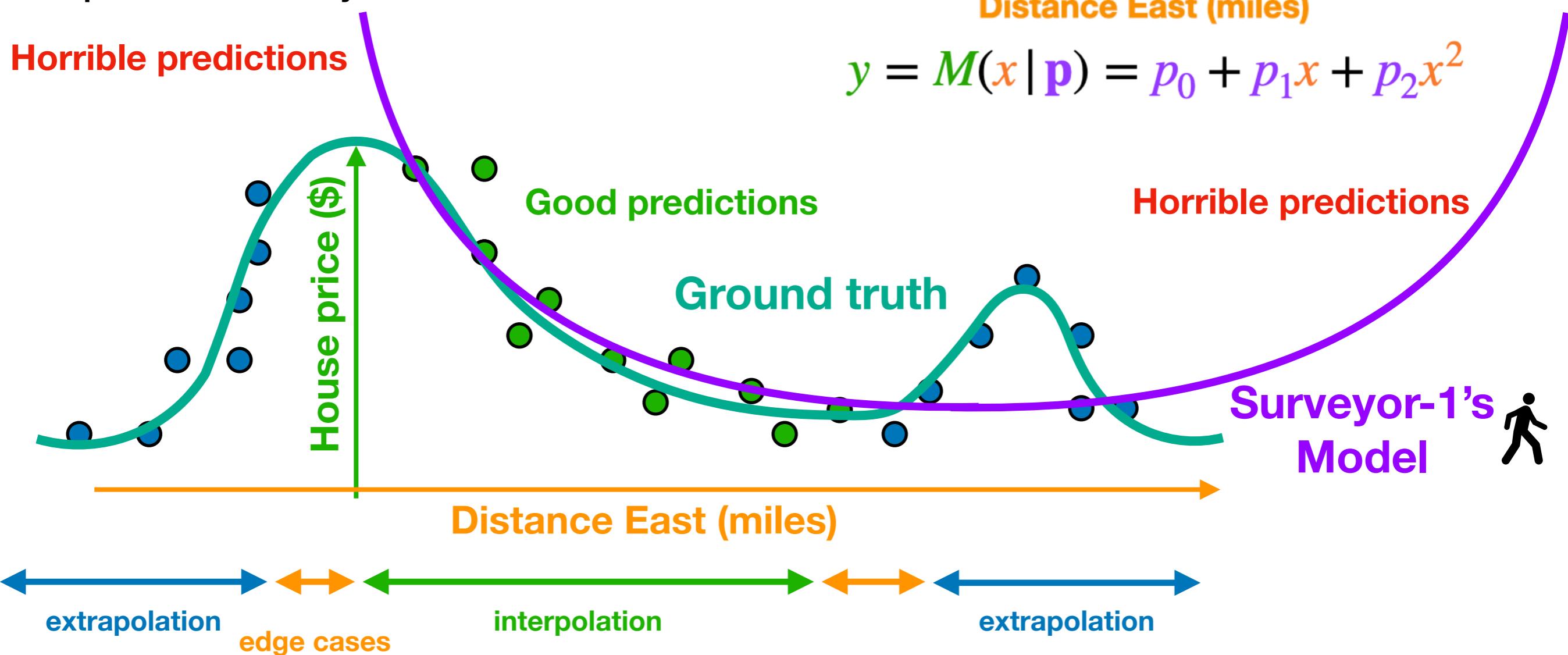
(AKA transferability)

A model's ability to extrapolate depends on;

- 1) Quality of the descriptor (inputs)
- 2) Amount and distribution data
- 3) Whether functional form of the model is similar to the ground truth

To improve transferability --> MORE DATA or better model

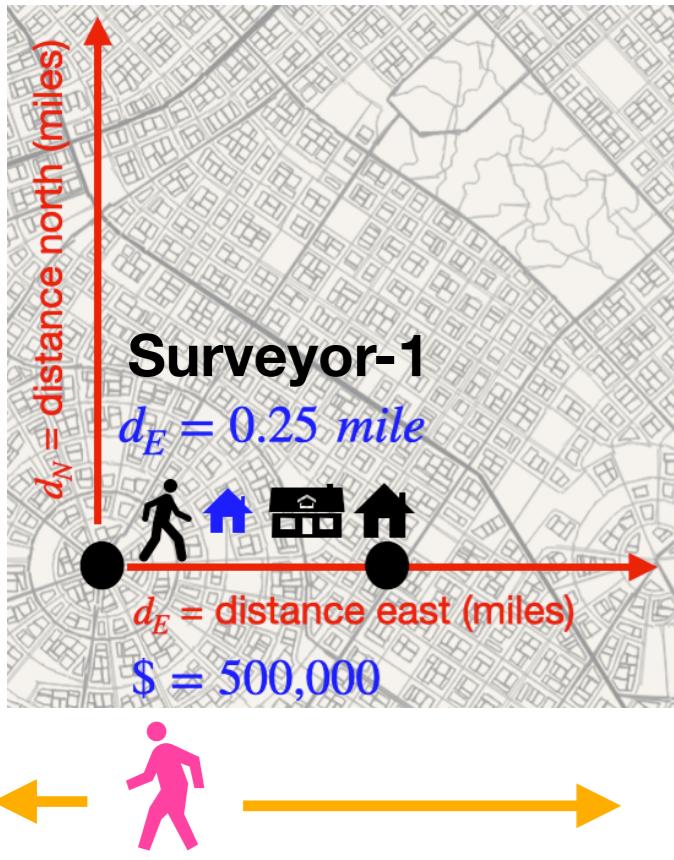
Horrible predictions



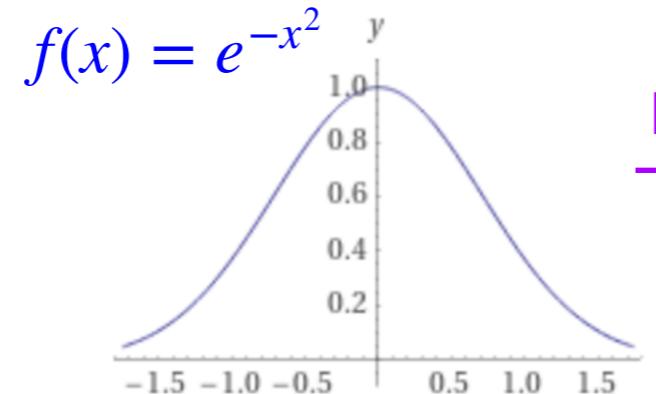
Interpolation: Predictions "Inside training region"

Most WELL TRAINED models can interpolate reasonably well

Surveyor-2 model

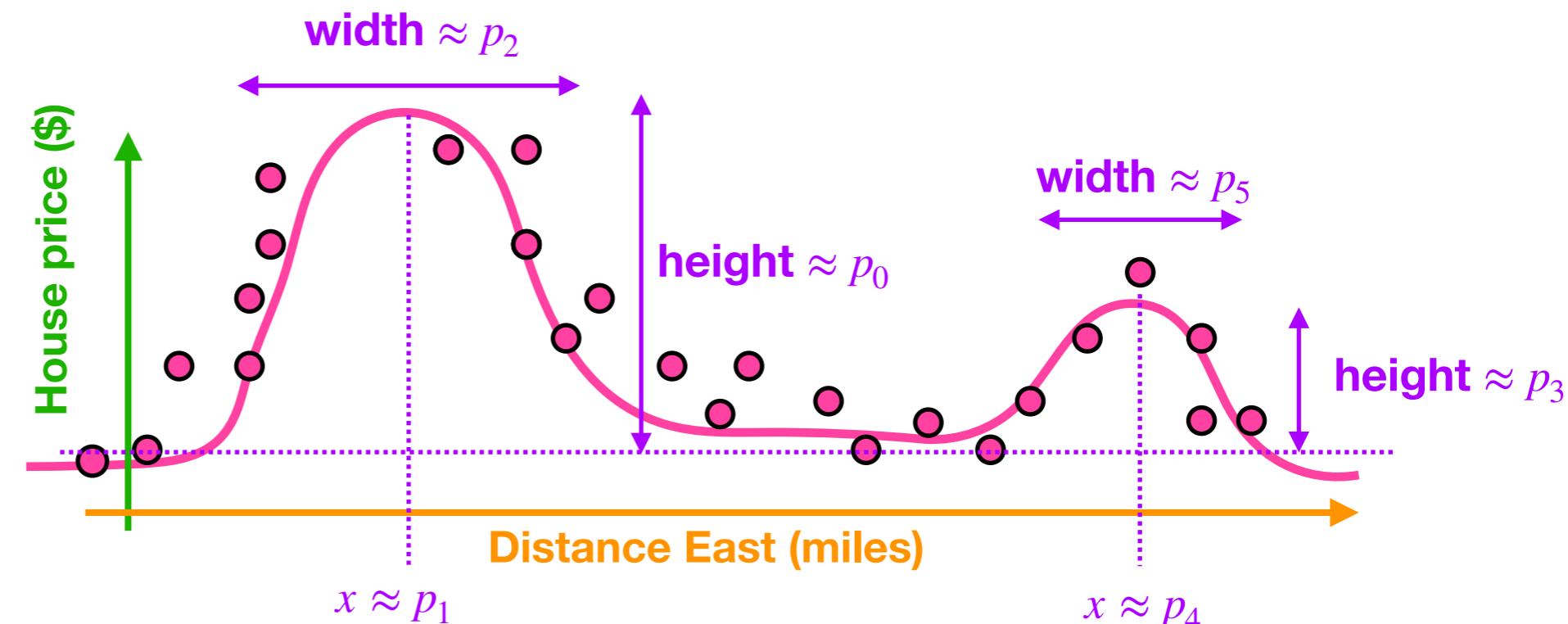


Parent function



Parent function

The purpose of the parameters is to stretch and bend the parent function to match "fit" the data



$M(x | p) = \text{Model} = \text{Functional form} + \text{Parameterization}$

$$M(x | p) = p_0 e^{-\left(\frac{x-p_1}{p_2}\right)^2} + p_3 e^{-\left(\frac{x-p_4}{p_5}\right)^2} + p_6$$

$$\mathbf{p} = (p_0, p_1, p_2, p_3, p_4, p_5, p_6)$$