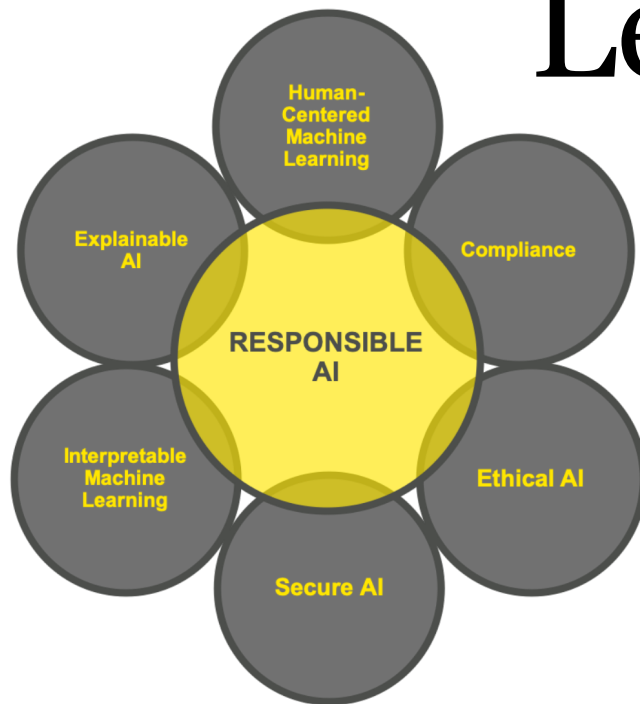


# Tufts COMP 135: Introduction to Machine Learning

<https://www.cs.tufts.edu/comp/135/2020f/>

# Responsible Machine Learning



## Model Cards for Model Reporting

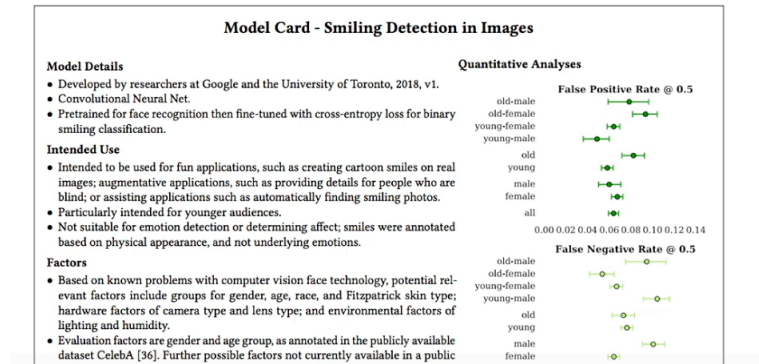


Image credit: Mitchell et al. '19

Image credit: h2o.ai

Prof. Mike Hughes

# Discussion Time!

See link to discussion guide on schedule

You have 15 minutes!

# Responsible ML: Takeaways

- Human-centered approach
  - Involve diverse stakeholders at every step
- Cyclical approach
  - Plan to Test, Release, Revise, Repeat
- Know your data: make a Datasheet
  - What is it capable of?
  - Limitations?
- Know your model: make a Model Card
  - What is it capable of?
  - Limitations? (e.g. most ML can't make causal claims)

# Review for Unit 5 Quiz

- SVMs for binary classification
  - Concepts: Support Vector, Hard Margin, Soft Margin
  - Hinge Loss
  - Compare/contrast with logistic regression
- Kernels
  - Definition of a kernel function
  - Examples: linear, squared-exponential, periodic
  - Primal vs dual view of prediction and training
  - Practical use: regression (like HW5), classification (with SVM)
- PCA
  - Encoding / decoding operations
  - Training objective
    - Minimize reconstruction error
  - Hyperparameter: How to select K?
- Fairness (at most 1 question)
  - How to tell if classifier is “fair”?
  - What metrics are appropriate?
    - Accuracy vs FPR vs TPR vs ....

# Two views of kernel prediction

1. Pick a kernel or feature transform (one implies other)

$$\phi(x_n) \in \mathbb{R}^G \quad k(x_i, x_j) = \phi(x_i)^T \phi(x_j)$$

2. Then choose a view below (both not always possible)

Primal (weights view, explicit feature vectors)

$$\theta \in \mathbb{R}^G \quad \hat{y}(x_i, \theta) = \theta^T \phi(x_i) = \sum_{g=1}^G \theta_g \cdot \phi(x_i)_g$$

Dual (kernel view, only inner prod of features)

$$\alpha \in \mathbb{R}^N \quad \hat{y}(x_i, \alpha, \{x_n\}_{n=1}^N) = \sum_{n=1}^N \alpha_n k(x_n, x_i)$$

# Why is kernel trick good idea?

Good with: squared exponential kernel  
Not so good with: linear kernel with raw features

Before (primal view)

Training problem seeks optimized vector of size  $G$

Prediction cost:

scales linearly with  $G$  (num. high-dim features)

requires  $G$  multiply ops plus  $G-1$  adds

After (dual view)

Training problem seeks optimized vector of size  $N$

Prediction cost:

scales linearly with  $N$  (num. train examples)

requires  $N$  evaluations of kernel plus  $N$  multiply/add

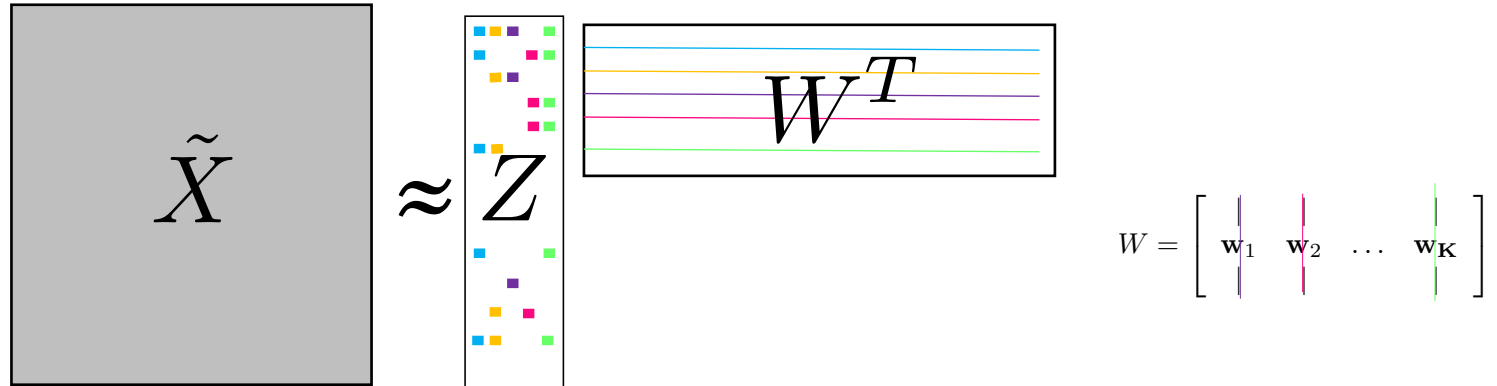
So we get some saving in runtime/storage with dual view if:

$G$  is bigger than  $N$

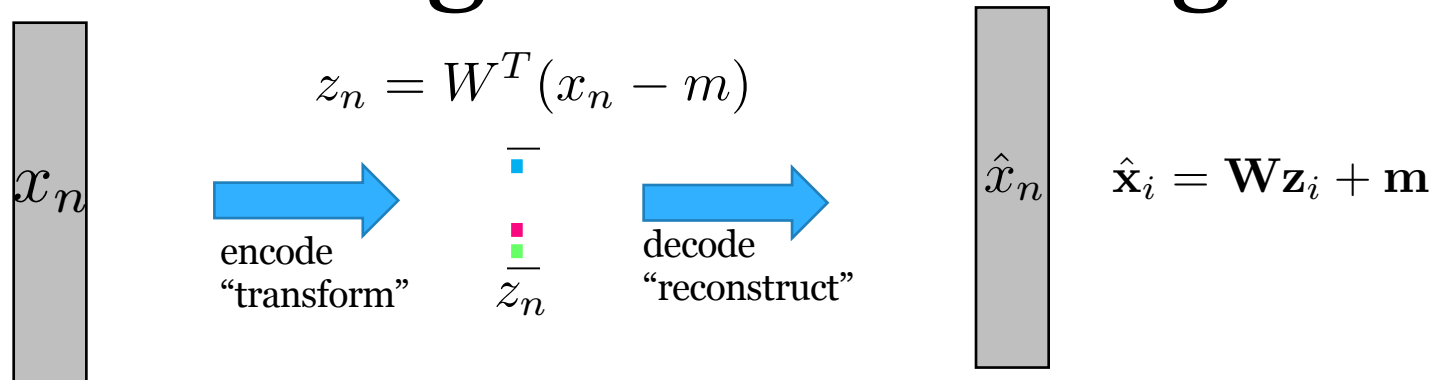
AND we can evaluate  $k$  fast (faster than a size  $G$  inner product)

Raw data:  $X$       Centered data:  $\tilde{X}$       Reconstructed data:  $\hat{X}$

# View: PCA as Matrix Factorization



# View: Encoding and Decoding



# Breakout: Practice Quiz Unit 5