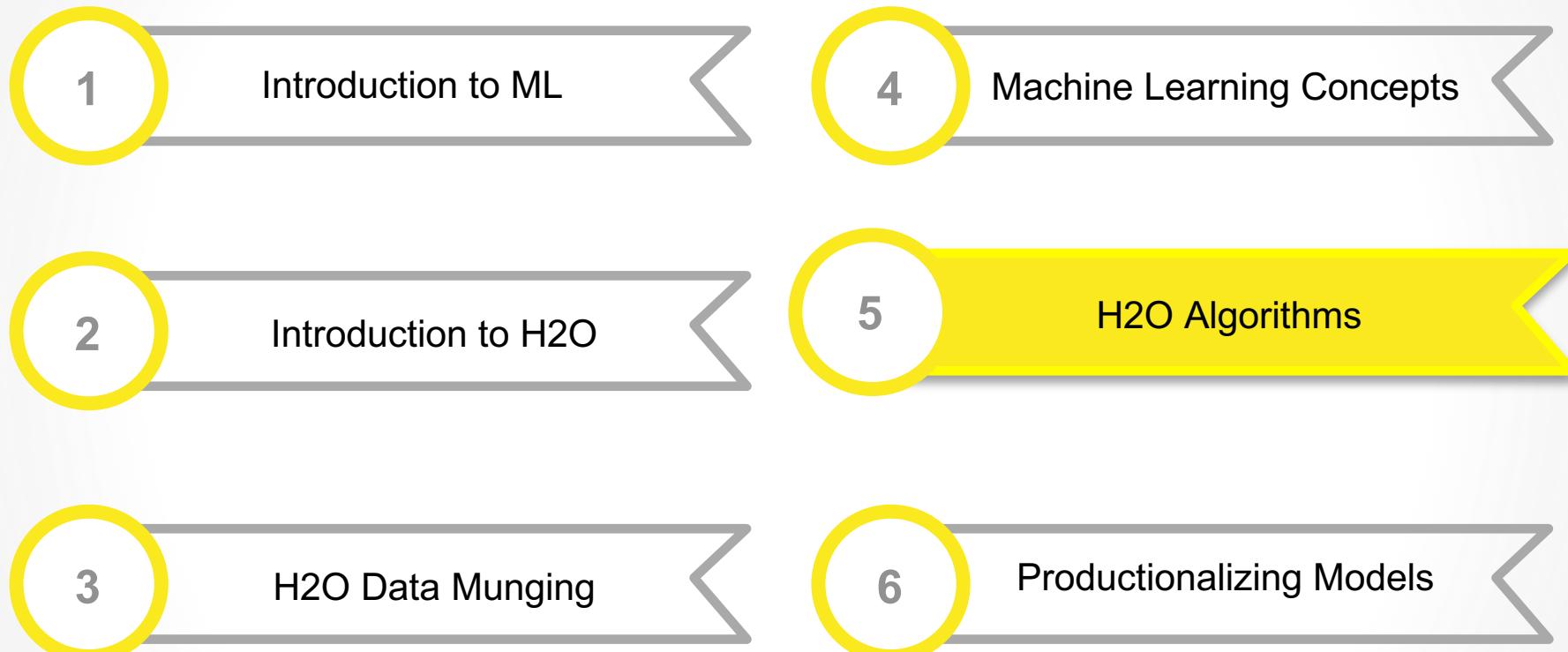


# H2O Training



# Supervised Learning in H2O

- **Regression and Classification Models** Exponential distributions including Poisson, Gamma, and Tweedie are available in addition to Bernoulli, Multinomial, and Gaussian distributions.
- **Fast and CPU Efficient** Parallel and distributed computation across multiple nodes and many cores.
- **Grid Search** Hyperparameter optimization allows the user to run through many parameters before selecting the best models.
- **Early Stopping** The user can specify the metric and the incremental change in the metric as convergence.
- **Stochastic** User can specify the sample rate that the algorithm will sample the column and row by for better generalization.
- **Model Output** The model is exportable as Java code and if you find the model overfitted after a certain number of trees, it is easy to reduce the number of trees in a POJO before putting it in production without rerunning model build.

Supervised Learning:

# **GENERALIZED LINEAR MODEL**

# What is the Generalized Linear Model (GLM)?

- Class of models that relate  $X$  (inputs) to  $Y$  (output)

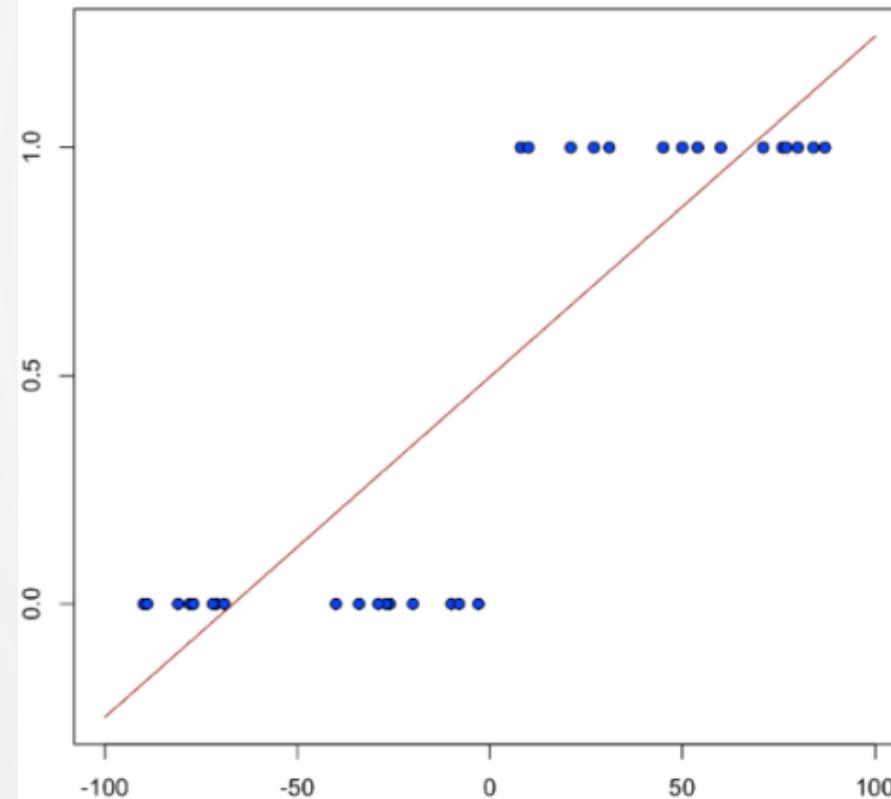
$$E(\mathbf{Y}) = \boldsymbol{\mu} = g^{-1}(\mathbf{X}\boldsymbol{\beta})$$

$$\text{Var}(\mathbf{Y}) = V(\boldsymbol{\mu}) = V(g^{-1}(\mathbf{X}\boldsymbol{\beta})).$$

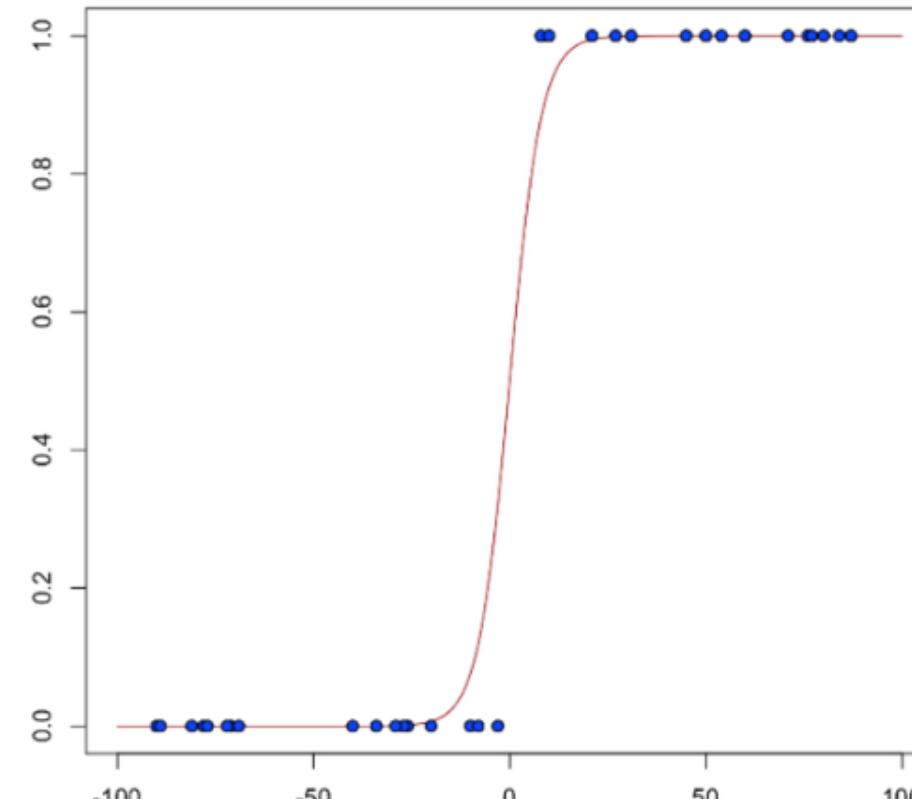
- Allows for a unification of models that have errors of the following form:
  - Normal (Gaussian)
  - Poisson
  - Gamma
  - Tweedie
  - Binomial (Logistic)
  - Multinomial
- MLE is found by iteratively reweighted least squares

# Same predictors, different family and link functions

Linear Regression fit  
(family=gaussian,link =identity)



Logistic Regression fit  
(family=binomial,link = logit)



# Pros and Cons of GLM's

## Pros

- Fast to fit
- Easy to interpret
- Well understood statistical properties
- Works for continuous, categorical, and count targets
- Transparent scoring functions

## Cons

- Dependent of external feature engineering
- Difficult to model non-linearities
- Difficult to model interaction effects
- Ill-defined with multicollinearity
  - Non-unique solutions

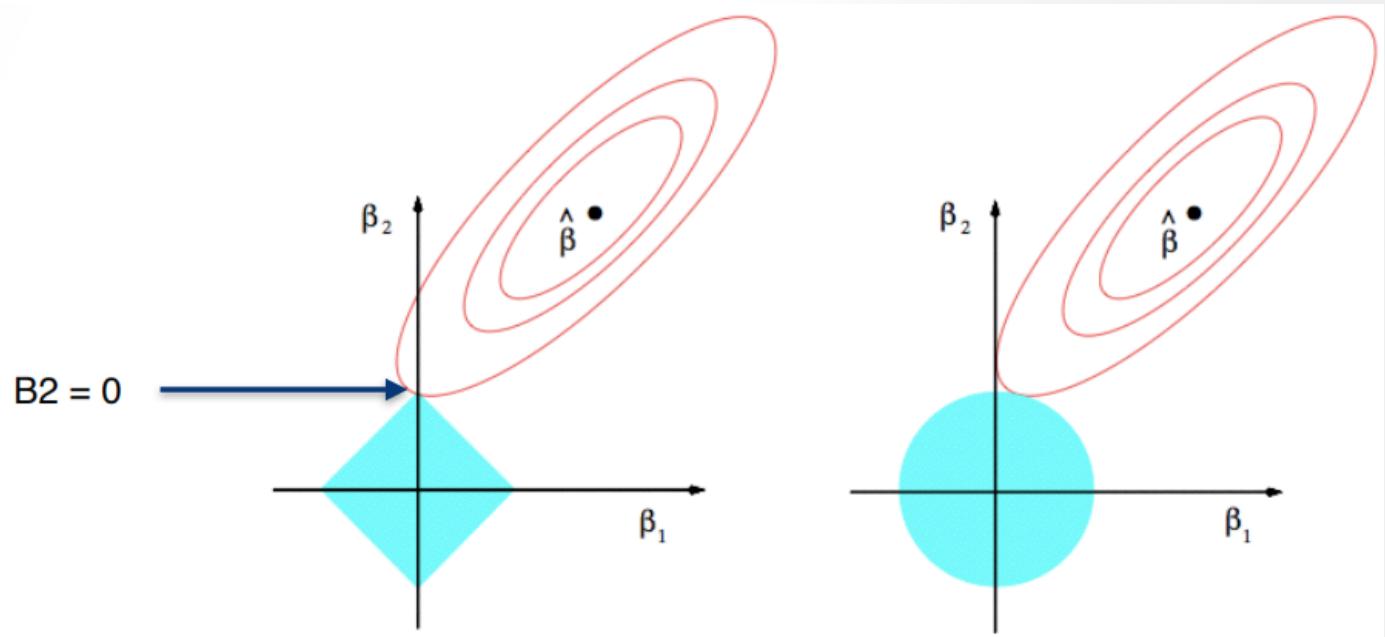
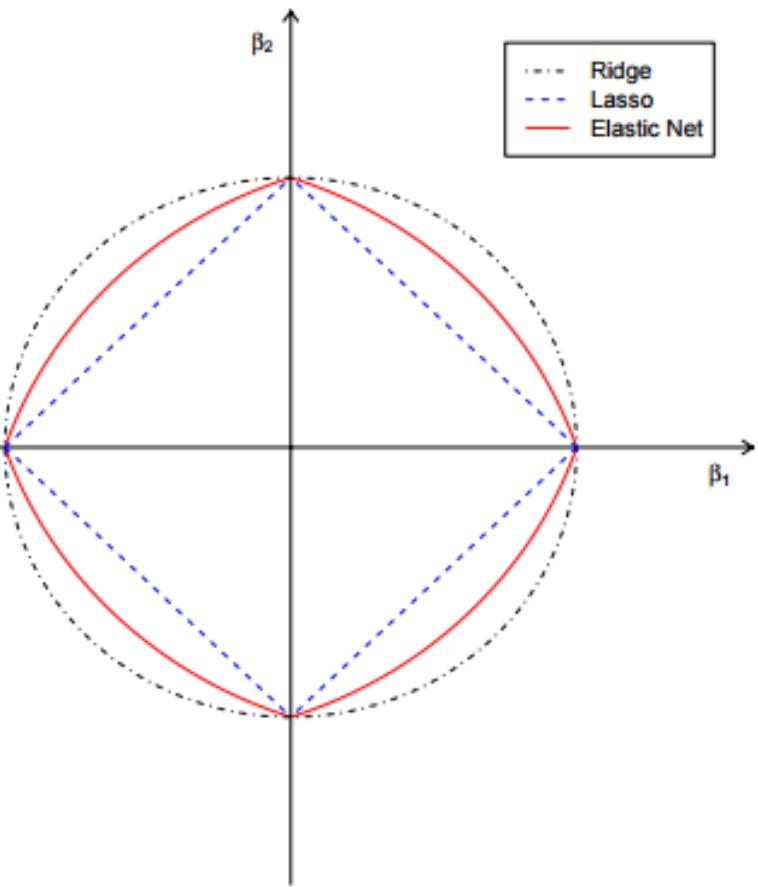
Supervised Learning:

# **PENALIZED GENERALIZED LINEAR MODEL(LASSO, RIDGE, ELASTIC NET)**

# Regularized Regression

- Want to build a parsimonious but interpretable model
- Shrink the number of predictors and/or size by imposing penalties on the estimated coefficients (L1, L2)
  - LASSO:  $\min_{\beta_0, \beta} \left\{ \frac{1}{N} \sum_{i=1}^N (y_i - \beta_0 - x_i^T \beta)^2 \right\}$  subject to  $\sum_{j=1}^p |\beta_j| \leq t$ .
    - picks one correlated variable, others discarded. Sparse.
  - Ridge:  $\text{minimize}_{\beta} \sum_{i=1}^n (y_i - \beta^T z_i)^2 \text{ s.t. } \sum_{j=1}^p \beta_j^2 \leq t$ 
    - correlated variables coefficients are pushed to the same value
  - Elastic Net:  $\arg \min_{\beta} \|\mathbf{y} - \mathbf{X}\beta\|^2 + \lambda_2 \|\beta\|^2 + \lambda_1 \|\beta\|_1$ 
    - sparse solution, correlated variables grouped, enter/ leave the model together

# Constraints for L1, L2 and L1&L2



LASSO Regression vs. Ridge Regression

# GLM Best Practices

- Regularization Selection
  - Explore a few values for alpha, e.g. 0, 0.25, 0.5, 0.75, 1
- Wide Data Sets (10K+ columns)
  - Iteratively Reweighted Least Squares (IRLS) fails with lambda = 0
    - IRLS requires  $p \times p$  Hessian matrix, where  $p = \#$  of coefficients
    - Could use Limited-memory BFGS (L-BFGS)
  - IRLS + lambda search works and is recommended
    - Use alpha  $\gg 0$
    - Can produce 1K+ non-zero coefficients
  - L-BFGS + L2 penalty works
  - L-BFGS + L1 penalty works, but may take a long time

# GLM Example

## scikit-learn-like interface for modeling

```
1 y = "income"
2 x = ["age", "workclass", "fnlwgt", "education", "marital-status", "occupation", "relationship",
3       "race", "sex", "capital-gain", "capital-loss", "hours-per-week", "native-country"]
```

```
1 from h2o.estimators.glm import H2OGeneralizedLinearEstimator
2 glm_0 = H2OGeneralizedLinearEstimator(family = "binomial", lambda_search = True,
3                                         nfolds = 5, seed = 123)
4 glm_0.train(x = x, y = y, training_frame = census_data, model_id = "income_glm_0")
```

glm Model Build progress: |██████████| 100%

```
1 from h2o.grid.grid_search import H2OGridSearch
2 glm_hyper_parameters = {"alpha": [0.5, 0.75, 1]}
3 glm_grid = H2OGridSearch(H2OGeneralizedLinearEstimator(family = "binomial", lambda_search = True,
4                                                               nfolds = 5, seed = 123),
5                           glm_hyper_parameters)
6 glm_grid.train(x = x, y = y, training_frame = census_data, grid_id = "income_glm_grid")
```

glm Grid Build progress: |██████████| 100%